

**An Interdisciplinary Approach to Human Pose  
Estimation: Application to Sign Language**

**Dissertation**

der Mathematisch-Naturwissenschaftlichen Fakultät  
der Eberhard Karls Universität Tübingen  
zur Erlangung des Grades eines  
Doktors der Naturwissenschaften  
(Dr. rer. nat.)

vorgelegt von  
Maria-Paola Forte  
aus Genua/Italien

Tübingen  
2025

Gedruckt mit Genehmigung der Mathematisch-Naturwissenschaftlichen Fakultät der Eberhard Karls Universität Tübingen.

Tag der mündlichen Qualifikation:	07.11.2025
Dekan:	Prof. Dr. Thilo Stehle
1. Berichterstatter:	Prof. Dr. Hendrik Lensch
2. Berichterstatterin:	Prof. Dr. Katherine J. Kuchenbecker
3. Berichterstatter:	Prof. Dr. Michael J. Black

---

This work is licensed under a CC BY-NC-ND 4.0 license:  
<https://creativecommons.org/licenses/by-nc-nd/4.0/legalcode.en>



# Abstract

Accessibility legislation, such as the Americans with Disabilities Act and the European Accessibility Act, mandates equal access to information and services for people with disabilities, including members of Deaf communities. While textual representations such as captions are commonly provided to meet these requirements, they present significant barriers for individuals whose primary language is sign language. Optimal accessibility would deliver information directly in sign language through videos of human interpreters. This approach is costly, time-consuming, and impractical for frequently updated or real-time content. AI-driven sign language avatars are a potential technological solution. Their training, however, requires high-quality motion data. This dissertation thus addresses the central research question: “Which technologies and methods can be used to collect high-quality (3D) data of sign language at scale?”

Vision-based human pose estimation offers a scalable approach. However, current methods encounter significant difficulties when applied to sign language due to its rapid hand movements and frequent self-occlusion and self-touch. This dissertation addresses these challenges by integrating two complementary approaches into pose-estimation methods: linguistic knowledge to improve hand pose accuracy and bioimpedance sensing to detect and enforce self-touch. Linguistic integration is achieved with the development of SGNify, a vision-based pose-estimation method that leverages universal linguistic rules from sign languages. SGNify formalizes linguistic constraints as computational priors to improve hand pose estimation and integrates facial expressions captured by specialized 3D face reconstruction models. The resulting 3D avatars quantitatively outperform those produced by standard methods, with proficient evaluators recognizing SGNify’s reconstructed signs at rates equivalent to the source videos.

Nevertheless, SGNify remains constrained by the inherent limitations of monocular videos; due to depth ambiguities and occlusions along the camera’s viewing axis, signs involving self-touch often result in inaccurate reconstructions biased toward non-contact configurations. To overcome this limitation, this dissertation explores the use of electrical bioimpedance to complement visual data during self-contact. Systematic measurements across diverse participants show that skin-to-skin contact causes significant bioimpedance magnitude changes at high frequencies (from approximately 240 kHz to 4.1 MHz), making bioimpedance a reliable measure for detecting self-touch across individuals.

These detected self-touch events are then used to resolve pose ambiguities through BioTUCH, a vision-based pose-estimation method complemented by bioimpedance sensing. BioTUCH applies a selective optimization strategy that refines the arm joints to cre-

ate physically plausible self-touch configurations. Quantitative evaluation demonstrates a substantial improvement in reconstruction accuracy compared to pure vision-based methods, enabling a more precise capture of self-contact in various scenarios, including sign language.

Through its contributions, this dissertation establishes novel methods for collecting high-quality sign language motion data at scale, addressing the data bottleneck impacting AI-driven sign language avatars. The techniques developed contribute to not only sign language technologies but also the broader field of human pose estimation.

# Zusammenfassung

Gesetzgebungen zur Barrierefreiheit wie das Americans with Disabilities Act und das European Accessibility Act schreiben einen gleichen Zugang zu Informationen und Dienstleistungen für Menschen mit Behinderungen vor, einschließlich Mitgliedern gehörloser Gemeinschaften. Während zur Erfüllung dieser Anforderungen üblicherweise textuelle Darstellungen wie Untertitel bereitgestellt werden, stellen diese erhebliche Barrieren für Personen dar, deren Primärsprache die Gebärdensprache ist. Optimale Barrierefreiheit hingegen würde die Informationen direkt in Gebärdensprache durch Videos menschlicher Dolmetscher vermitteln. Ein solcher Ansatz ist jedoch kostspielig, zeitaufwendig und unpraktisch für häufig aktualisierte oder Echtzeit-Inhalte. KI-gesteuerte Gebärdensprach-Avatare sind eine potenzielle technologische Lösung. Ihr Training erfordert jedoch hochwertige Bewegungsdaten. Diese Dissertation befasst sich daher mit der zentralen Forschungsfrage: „Welche Technologien und Methoden können verwendet werden, um hochwertige (3D-)Daten von Gebärdensprache in großem Maßstab zu erfassen?“

Vision-basierte menschliche Posenschätzung bietet einen skalierbaren Ansatz. Aktuelle Methoden stoßen jedoch bei der Anwendung auf Gebärdensprache auf erhebliche Schwierigkeiten aufgrund von schnellen Handbewegungen und häufiger Selbstverdeckung und Selbstberührung. Diese Dissertation adressiert diese Herausforderungen durch die Integration zweier komplementärer Ansätze in der Posenschätzung: Einerseits wird linguistisches Wissen zur Verbesserung der Handposengenauigkeit genutzt, andererseits kommt Bioimpedanz-Sensorik zum Einsatz, um Selbstberührung zu erkennen und durchzusetzen. Die linguistische Integration wird durch die Entwicklung von SGNify erreicht, einer vision-basierten Methode zur Posenschätzung, die universelle linguistische Regeln aus Gebärdensprachen nutzt. SGNify formalisiert linguistische Einschränkungen als computergestützte Vorannahmen zur Verbesserung der Handposenschätzung und integriert Gesichtsausdrücke, die durch spezialisierte 3D-Gesichtsrekonstruktionsmodelle erfasst werden. Die resultierenden 3D-Avatare übertreffen quantitativ diejenigen, die durch Standardmethoden erzeugt werden, wobei versierte Begutachter rekonstruierte Gebärden von SGNify mit Raten erkennen, die denen der Originalvideos entsprechen.

Dennoch bleibt SGNify durch die inhärenten Einschränkungen monokularer Videos begrenzt; aufgrund von Tiefenambiguitäten und Verdeckungen entlang der Kamera-Sichtachse führen Bewegungen mit Selbstberührung oft zu ungenauen Rekonstruktionen, die zu kontaktlosen Konfigurationen verzerrt sind. Um diese Einschränkung zu überwinden, untersucht diese Dissertation die Verwendung elektrischer Bioimpedanz zur Ergänzung visueller Daten während Selbstkontakt besteht. Systematische Messungen über diverse

Teilnehmer hinweg zeigen, dass Haut-zu-Haut-Kontakt signifikante Änderungen der Bioimpedanzmagnitude bei hohen Frequenzen (von ungefähr 240 kHz bis 4,1 MHz) verursacht, was Bioimpedanz zu einem zuverlässigen Maß für die Erkennung von Selbstberührung über Individuen hinweg macht.

Diese erkannten Selbstberührungseignisse werden dann verwendet, um Posenambiguitäten durch BioTOUCH—eine vision-basierte Posenschätzungsmethode, die durch Bioimpedanz-Sensorik ergänzt wird—aufzulösen. BioTOUCH wendet eine selektive Optimierungsstrategie an, die die Armgelenke verfeinert, um physikalisch plausible Selbstkontakt-Konfigurationen zu erzeugen. Eine quantitative Evaluation demonstriert eine substantielle Verbesserung der Rekonstruktionsgenauigkeit im Vergleich zu rein vision-basierten Methoden und ermöglicht eine präzisere Erfassung von Selbstkontakt in verschiedenen Szenarien, einschließlich Gebärdensprache.

Durch ihre Beiträge etabliert diese Dissertation neuartige Methoden zur Erfassung hochwertiger Gebärdensprach-Bewegungsdaten in großem Maßstab und adressiert den Datenengpass, der KI-gesteuerte Gebärdensprach-Avatare beeinträchtigt. Die entwickelten Techniken tragen nicht nur zu Gebärdensprachtechnologien bei, sondern auch zum breiteren Feld der menschlichen Posenschätzung.

# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
<b>2</b>	<b>Background</b>	<b>5</b>
2.1	Linguistic Properties of Sign Language . . . . .	5
2.1.1	Phonological Parameters of Sign Languages . . . . .	5
2.1.2	Classification of Signs Based on Phonological Properties . . . . .	6
2.1.3	Self-Touch in Sign Language . . . . .	6
2.2	Human Pose Estimation for Sign Language . . . . .	8
2.2.1	Human Body Models . . . . .	8
2.2.2	Optimization-Based Human Pose Estimation . . . . .	9
2.2.3	Challenges for Sign Language Capture . . . . .	10
2.2.4	Current Methods for Sign Language Capture . . . . .	11
2.3	Sensing Self-Contact . . . . .	11
2.3.1	Electrical Bioimpedance Sensing . . . . .	12
2.4	Integrating Linguistics and Bioimpedance Sensing into Human Pose Estimation . . . . .	14
<b>3</b>	<b>Integrating Linguistic Priors into Human Pose Estimation</b>	<b>17</b>
3.1	SGNify . . . . .	19
3.1.1	SMPLify-SL: Baseline for Sign Language Video . . . . .	19
3.1.2	Linguistic Constraints . . . . .	20
3.1.3	Automatization . . . . .	23
3.2	Dataset . . . . .	26
3.3	Experiments . . . . .	27
3.3.1	Quantitative Evaluation . . . . .	27
3.3.2	Qualitative Results . . . . .	29
3.3.3	Perceptual Study . . . . .	29
3.3.4	SGNify Extensions . . . . .	35
3.4	Discussion . . . . .	37
3.4.1	Limitations and Future Work . . . . .	38
<b>4</b>	<b>Detecting Self-Touch Using Bioimpedance</b>	<b>39</b>
4.1	Materials and Methods . . . . .	41
4.1.1	Experimental Setup . . . . .	41
4.1.2	Experimental Protocol . . . . .	42

4.1.3	Dataset . . . . .	44
4.1.4	Statistical Analysis . . . . .	45
4.2	Results . . . . .	48
4.2.1	Detection of Self-Touch Poses . . . . .	48
4.2.2	Influence of Body Parts and Skin Contact Area on Skin-to-Skin Bioimpedance Changes . . . . .	52
4.3	Discussion . . . . .	53
4.3.1	Limitations and Future Work . . . . .	55
<b>5</b>	<b>Integrating Bioimpedance Sensing into Human Pose Estimation</b>	<b>57</b>
5.1	Method . . . . .	59
5.1.1	Bioimpedance Sensing . . . . .	60
5.1.2	BioTOUCH . . . . .	60
5.2	Dataset . . . . .	63
5.3	Experiments . . . . .	65
5.3.1	Evaluation of Self-Contact Detection . . . . .	65
5.3.2	Quantitative Evaluation . . . . .	65
5.3.3	Qualitative Evaluation . . . . .	68
5.3.4	In-the-Wild Capture . . . . .	69
5.3.5	Application to Sign Language . . . . .	70
5.4	Discussion . . . . .	70
5.4.1	Limitations and Future Work . . . . .	71
<b>6</b>	<b>Discussion</b>	<b>73</b>
6.1	Summary of Contributions . . . . .	73
6.2	Implications and Applications . . . . .	74
6.3	Limitations and Future Work . . . . .	75
6.4	Conclusion . . . . .	76
	<b>Bibliography</b>	<b>77</b>

# Chapter 1

## Introduction

“Thanks” is one of the first words that people learn when studying a new language. In American Sign Language (ASL), this sign starts with the dominant hand open and palm facing the lips, fingers extended and held together touching the mouth. The hand then moves forward and a bit down in the direction of the person being thanked, keeping a smile. Processing this textual description may result in physical mimicry, mental visualization, or the recognition that a textual description is inadequate. In contrast, an annotated photograph like the one shown in Fig. 1.1 would greatly facilitate understanding. The effectiveness of visual representation over textual description highlights the defining characteristic of sign language (SL): its visual-spatial modality. SL is perceived exclusively visually [72], and each sign combines simultaneous articulations of the hands, face, and upper body, utilizing the 3D space around the signer.



Figure 1.1: Visual representation of the sign THANKS in ASL. The arrow indicates the forward and slightly downward movement of the hand.

SL serves as the primary communication method for the Deaf community; with disabling hearing loss projected to affect one in every ten people by 2050 [92], ensuring accessible communication technologies becomes increasingly critical. To ensure timely access to information, AI-driven avatars that translate written text into SL videos could play a pivotal role. However, creating such avatars requires high-quality training data that captures the nuanced visual-spatial properties of SL, including precise hand configurations, facial expressions, and self-touch<sup>1</sup> events.

Current approaches to capturing this data encounter significant challenges. Vision-based methods, whether using 2D or 3D representations [46], struggle with rapid hand movements that cause motion blur. While monocular cameras offer scalable capture, they are limited by self-occlusions and depth ambiguity, making it difficult to reliably determine whether a hand is touching the face or simply positioned near it. Furthermore, although 3D meshes offer advantages over other representations, such as view-invariant features, detailed hand reconstruction [73], and robust co-articulation modeling [98], they still fail to accurately capture even simple signs like THANKS<sup>2</sup>.

This persistent failure reveals a critical research gap: despite decades of linguistic research establishing the systematic structure of signs [4] and emerging sensing technologies that could detect contact events, no existing approach integrates either of these complementary sources of information to overcome the inherent limitations of vision-only systems. This dissertation thus enables more accurate 3D data collection for SL, a process we define as Sign Language Capture (SLC), by introducing novel methods that combine linguistic knowledge and bioimpedance sensing with vision-based human pose estimation (HPE); linguistic knowledge guides anatomically and linguistically correct pose estimation, while bioimpedance resolves contact ambiguities that monocular vision cannot disambiguate.

The interdisciplinary approach to HPE developed in this dissertation is tailored to the unique challenges of SL but also has broader applications in domains requiring precise reconstruction of self-contact, such as healthcare monitoring, sports performance analysis, and extended reality applications. Throughout this work, ethical considerations have been central, including collaborations with Deaf community members to ensure that technological developments accurately represent the culture.

The next chapter (**Chapter 2**) presents the background information needed to understand the subsequent chapters. **Chapters 3, 4, and 5** form the core of this dissertation, addressing key challenges of SLC. The dissertation concludes with an overall discussion in **Chapter 6**, which reflects on the main contributions of this work and, considering its

---

<sup>1</sup>Throughout this dissertation, we use “self-touch” when referring to hand-initiated contact (*e.g.*, hand touching face or hands clasped together). We use the broader term “self-contact” for any contact between any body parts (*e.g.*, ankles touching or elbow to torso). Both self-touch and self-contact can happen either directly on skin or through clothing.

<sup>2</sup>A gloss, such as THANKS, is a written representation of a sign, typically using words from the dominant spoken language, *e.g.*, English for ASL or German for German Sign Language (DGS), that function as a label for the sign. By convention, glosses are written with small capital letters.

---

limitations, outlines directions for future research. The detailed outline of **Chapters 3, 4, and 5** is as follows:

**Chapter 3: Integrating Linguistic Priors into Human Pose Estimation.** This chapter presents SGNify, a method for reconstructing expressive 3D signing avatars from monocular SL videos. The chapter details how linguistic knowledge from SL phonology can be formalized into mathematical constraints to improve hand reconstruction accuracy. Specifically, we develop and implement two key linguistic constraints: hand-pose symmetry and hand-pose invariance, which are derived from universal principles applicable across different SLs. The chapter explains the system’s pipeline, including the sign-group classifier that automatically determines which linguistic constraints to apply based on the observed signing. The evaluation section presents both quantitative results, comparing SGNify with previous HPE methods on a newly collected motion-capture dataset, and qualitative analysis through a perceptual study with proficient signers. Finally, the chapter discusses the implications of this work for SL learning and accessibility technology. This chapter was previously published at CVPR 2023 [24].

**Chapter 4: Detecting Self-Touch Using Bioimpedance.** This chapter explores the use of electrical bioimpedance for detecting self-touch poses. It describes the experimental setup, including electrode placement at the wrists and signal acquisition. It then presents a systematic investigation of 27 genuine self-touch poses and six adversarial mid-air gestures across 30 diverse participants. The results section provides a comprehensive analysis of how various self-touch poses affect bioimpedance measurements, examining the difference between skin-to-skin and skin-to-clothing poses, identifying the optimal frequency ranges for detection, and quantifying the influence of factors such as body parts involved, skin contact area, and individual characteristics like sex, ethnicity, and body mass index (BMI). The chapter also discusses the practical implications of these findings for wearable technology, including design considerations for non-invasive self-touch detection systems. This chapter was previously published in the IEEE Transactions on Instrumentation and Measurement [28].

**Chapter 5: Integrating Bioimpedance Sensing into Human Pose Estimation.** This chapter integrates the findings from Chapter 4 into vision-based HPE by introducing BioTUCH. The chapter details the BioTUCH pipeline, including the self-contact detection algorithm and the optimization strategy that refines arm pose estimates from off-the-shelf HPE methods. It then presents a newly collected dataset of synchronized RGB videos, bioimpedance measurements, and 3D motion capture data, comprising 82 common self-touch gestures and nine adversarial non-contact gestures performed by three participants in different postures. The evaluation section provides quantitative results demonstrating an average relative improvement of 11.67% in self-touch reconstruction accuracy, along with an absolute increase of 31.60 percentage points in correctly reconstructing contact when it is present in the ground truth, compared to the HPE methods used as input. The chapter also describes the development and testing of a miniaturized wearable bioimpedance-based sensor that enables practical capture of self-contact at scale. Finally, it discusses the broader implications of this multimodal approach for

creating high-quality pseudo-ground truth data that can improve the training of HPE methods, benefiting applications ranging from SL to human behavior analysis. This chapter was previously published at ICCV 2025 [25].

# Chapter 2

## Background

This chapter provides the interdisciplinary background necessary to understand the challenges of SLC and how this dissertation addresses them. While SLC appears to be a computer-vision problem, its effective solution requires understanding both the linguistic structure of SLs and the fundamental limitations of current HPE methods.

### 2.1 Linguistic Properties of Sign Language

The systematic linguistic study of SLs began with William Stokoe's groundbreaking work in the 1960s, which demonstrated that ASL exhibited the fundamental properties of a natural language with its own grammar and structure [83]. This work shifted the paradigm from viewing SLs as mere gestural systems to recognizing them as natural languages.

During the 1970s and 1980s, researchers including Klima and Bellugi further developed the phonological analysis of SLs, establishing that signs consist of discrete components analogous to phonemes in spoken languages [44]. Battison's work in particular provided foundational insights into the constraints that govern how signs can be formed, establishing a typology that remains influential in SL research [4].

#### 2.1.1 Phonological Parameters of Sign Languages

Each sign is composed of a combination of the following five discrete parameters:

- **Handshape:** the configuration of the fingers. SLs have a constrained inventory of handshapes, with each SL utilizing a subset of all possible finger configurations.
- **Location:** the spatial position where the sign is articulated. It may involve parts of the body (*e.g.*, forehead or torso) or exist in the neutral signing space in front of the signer.
- **Movement:** the dynamic component of a sign. It includes the direction, speed, repetition, and manner (*e.g.*, straight, circular, or zigzag).

- **Orientation:** the direction in which the palm and fingers face during the sign.
- **Non-manual features:** facial expressions, mouthing, head tilts, shoulder movements, and eye gaze. These features are essential for conveying syntactic structures (e.g., questions or negations) and modifying lexical meaning.

These parameters combine simultaneously rather than sequentially, which is a crucial distinction from the linear arrangement of phonemes in spoken languages. This simultaneity creates a dense information packaging system that efficiently utilizes the 3D space [79].

### 2.1.2 Classification of Signs Based on Phonological Properties

Battison introduced a typology that divides signs based on distinct types of motor acts [4]:

- **Type 0:** one-handed signs produced with the dominant hand, articulated in free space without contact with any body parts.
- **Type X:** one-handed signs produced with the dominant hand, contacting the body at any location except the opposite hand.
- **Type 1:** two-handed signs in which both hands have identical handshapes and perform symmetrical or alternating movements. These signs typically conform to the *Symmetry Condition*<sup>3</sup>.
- **Type 2:** two-handed signs in which the dominant hand moves while the non-dominant hand acts as a stationary base; both hands share the same handshape.
- **Type 3:** The dominant hand moves while the non-dominant hand acts as a stationary base and has a different, typically neutral, handshape. These signs generally conform to the *Dominance Condition*<sup>4</sup>.

In addition to these five basic categories, Battison also proposed a sixth type, **Type C**, to account for *compound signs*, which combine two or more of the above types within a single lexical item [4].

### 2.1.3 Self-Touch in Sign Language

Self-touch is a phonologically distinctive feature in SLs, justifying Battison's separate classification of Type 0 and Type X signs. However, its frequency is not documented. To

---

<sup>3</sup>Symmetry Condition: if both hands move independently in a two-handed sign, they must have identical handshapes, locations, and types of movement, with the movement being either synchronous or alternating in a 180° phase relationship.

<sup>4</sup>Dominance Condition: in two-handed signs with different handshapes, the non-dominant hand must remain stationary, and its handshape must come from a limited set of unmarked forms (typically B, A, S, 1, C, O, or 5 in ASL).

quantify contact frequency, we can leverage HamNoSys [33] annotations. HamNoSys (Hamburg Sign Language Notation System) is a universal SL phonetic transcription system that can be used to represent all hand poses and movements that constitute a sign; *i.e.*, a HamNoSys annotation contains sufficient information to reproduce the represented sign. We constructed an Extended Backus–Naur Form (EBNF) grammar (see Fig. 2.1) to parse HamNoSys annotations and extract labels to compute statistics on the signs, providing a more robust framework than previous decision tree methods [58].

Parsing the HamNoSys entries of ASL and Polish Sign Language (PJM), we computed that over 50% of their signs involve contact. Furthermore, its presence can differentiate between minimal pairs (*i.e.*, pair of signs that differ by only one parameter); in ASL, DRY and UGLY differ solely by the contact of the hand with the chin in DRY. Finally, contact may occur between hands (*e.g.*, NAME) or with other upper body parts like the face (*e.g.*, DRY or THANKS) or the torso (*e.g.*, PLEASE), and can be realized as quick taps (*e.g.*, THANKS), repeated contacts (*e.g.*, NAME), or sustained holds (*e.g.*, PLEASE).

```

⟨hns⟩ ::= [SYMMETRY] ⟨block⟩

⟨block⟩ ::= [⟨handshape_block⟩ | ⟨non_handshape_block⟩]*

⟨handshape_block⟩ ::= HANDSHAPE [HANDSHAPE_MODIFIER | HANDSHAPE_FINGER_LOCATION]*

⟨non_handshape_block⟩ ::= ⟨par⟩ | ⟨seq⟩ | ⟨fusion⟩ | EXTENDED_FINGER_LOCATION | PALM_ORIENTATION | MOVEMENT | MOVEMENT_MODIFIER | LOCATION | LOCATION_MODIFIER | OTHER_SYMBOL_NO_GROUP

⟨par⟩ : HAMPARBEGIN ⟨block⟩ [HAMPLUS ⟨block⟩] HAMPAREND

⟨seq⟩ : HAMSEQBEGIN ⟨block⟩ HAMSEQEND

⟨fusion⟩ : HAMFUSIONBEGIN ⟨block⟩ HAMFUSIONEND

```

Figure 2.1: Constructed HamNoSys EBNF grammar showing the hierarchical structure of HamNoSys notation. Blocks can contain handshape information or other elements like location and movement, with support for parallel, sequential, and fusion operations.

## 2.2 Human Pose Estimation for Sign Language

The development of accurate HPE has been a gradual process. Human body modeling has evolved significantly since the early attempts in the 1970s, which typically represented the human bodies as simple stick figures. These early models eventually gave way to more sophisticated approaches that could represent the body surface, beginning with volume-based models using cylinders, ellipsoids, or spheres. A major advancement came in the late 1980s with the development of models where skin deformation was driven by an underlying skeleton, allowing for more realistic representation of human movement. However, these early approaches could not capture the complexity of actual human anatomy and movement [66].

The field was transformed by the introduction of data-driven approaches, beginning with SCAPE [1] in 2005, which learned how the human body deforms with pose from real-world scan data. This advance was followed by SMPL [55], which combined an underlying skeletal structure with linear blend skinning (*i.e.*, smoothly blending the effects of multiple joint transformations to deform the body mesh), learned from large datasets, making it particularly suitable for animation and computer vision applications.

The development of specialized models for different body parts followed, with MANO for hands [77], FLAME for faces [54], and later models such as SMPL-X [71], which incorporate articulated fingers and facial expressions alongside full-body pose representation. These integrated models are particularly relevant for SL, where handshape, facial expression, and upper body posture all carry linguistic information.

### 2.2.1 Human Body Models

#### SMPL

The SMPL model [55] provides a mathematical foundation for modern human body modeling and can be formulated as:

$$M(\beta, \theta; \Phi) : \mathbb{R}^{|\beta| \times |\theta|} \rightarrow \mathbb{R}^{3N} \quad (2.1)$$

where  $M$  maps body shape ( $\beta \in \mathbb{R}^{|\beta|}$ ) and body pose ( $\theta \in \mathbb{R}^{|\theta|}$ ) via learned parameters  $\Phi = \{\bar{\mathbf{T}}, \mathcal{J}, \mathcal{W}, \mathcal{S}, \mathcal{P}\}$  to a 3D mesh with  $N = 6,890$  vertices ( $\mathbf{V} \in \mathbb{R}^{3N}$ ), connected through triangular faces, and 23 joints, connected through a skeletal rig (with an additional joint for global orientation). The learned parameters ( $\Phi$ ) include:

- $\bar{\mathbf{T}}$ : the template mesh in rest pose (a neutral pose with arms extended).
- $\mathcal{J}$ : a function that computes joint locations from the rest pose mesh.
- $\mathcal{W}$ : blend weights that determine how much each vertex is affected by joint rotations.

- $\mathcal{S}$ : shape blend shapes that capture how body shape varies across individuals.
- $\mathcal{P}$ : pose blend shapes that model pose-dependent deformations.

The SMPL function is thus defined as:

$$M(\beta, \theta) = W(T_P(\beta, \theta), \mathbf{J}(\beta), \theta, \mathcal{W}), \quad (2.2)$$

where

$$T_P(\beta, \theta) = \bar{\mathbf{T}} + B_S(\beta; \mathcal{S}) + B_P(\theta; \mathcal{P}), \quad (2.3)$$

represents the posed template with shape- and pose-dependent vertex displacements applied via  $B_S(\beta; \mathcal{S})$  and  $B_P(\theta; \mathcal{P})$ , respectively.

### SMPL-X

SMPL-X [71] is the expressive extension of SMPL. It adds 15 joints per hand and the facial expression parameters ( $\psi$ ):

$$M(\beta, \theta, \psi; \Phi) : \mathbb{R}^{|\beta| \times |\theta| \times |\psi|} \rightarrow \mathbb{R}^{3N}, \quad (2.4)$$

$$M(\beta, \theta, \psi) = W(T_P(\beta, \theta, \psi), \mathbf{J}(\beta), \theta, \mathcal{W}), \quad (2.5)$$

where

$$T_P(\beta, \theta, \psi) = \bar{\mathbf{T}} + B_S(\beta; \mathcal{S}) + B_P(\theta; \mathcal{P}) + B_E(\psi; \mathcal{E}), \quad (2.6)$$

and  $B_E(\psi; \mathcal{E})$  represents blend shapes for facial expressions, with  $\mathcal{E}$  being the learned facial expression blend shapes. The extended parameter set is thus  $\Phi = \{\bar{\mathbf{T}}, \mathcal{J}, \mathcal{W}, \mathcal{S}, \mathcal{P}, \mathcal{E}\}$ . As such, SMPL-X features a more detailed mesh topology with  $N = 10,475$  vertices.

## 2.2.2 Optimization-Based Human Pose Estimation

There are two common approaches to estimating a person’s body pose and shape in 3D from images: the parameters of a human body model like SMPL-X can be either optimized or regressed. During optimization, the primary objective is to minimize the Euclidean distance between the estimated 3D joints projected onto the image and the ground-truth (GT) 2D joints. Pose and shape priors prevent unrealistic estimates. Regression networks are instead usually trained on datasets of images with paired 3D GT data, *e.g.*, from motion capture. In general, optimization-based methods are more computationally intensive but produce more accurate results when limited training data is available. Due to the scarcity of GT data for the applications presented in this dissertation, we build on SMPLify-X [71], an optimization-based approach for fitting SMPL-X bodies to the 2D joint locations estimated from keypoints detectors.

SMPLify-X optimizes pose parameters ( $\theta$ ), shape parameters ( $\beta$ ), and facial expression parameters ( $\psi$ ) to minimize the re-projection error between estimated 2D joints

( $J_{est}$ ) and projected 3D model joints ( $R_\theta(\mathbf{J}(\boldsymbol{\beta}))$ ), along with multiple regularization terms. The optimization objective is:

$$E(\boldsymbol{\beta}, \boldsymbol{\theta}, \boldsymbol{\psi}) = E_J + \lambda_{\theta_b} E_{\theta_b} + \lambda_{\theta_f} E_{\theta_f} + \lambda_{m_h} E_{m_h} + \lambda_\alpha E_\alpha + \lambda_\beta E_\beta + \lambda_\epsilon E_\epsilon + \lambda_C E_C, \quad (2.7)$$

where  $E_J$  represents the joint re-projection loss computed as:

$$E_J(\boldsymbol{\beta}, \boldsymbol{\theta}, K, J_{est}) = \sum_{\text{joint } i} \gamma_i \omega_i \rho(\Pi_K(R_\theta(\mathbf{J}(\boldsymbol{\beta}))_i) - J_{est,i}), \quad (2.8)$$

with  $\rho$  being a robust Geman-McClure penalty function [29],  $\omega_i$  the detection confidence,  $\gamma_i$  per-joint weights for annealed optimization, and  $\Pi_K$  the camera projection with intrinsic parameters ( $K$ ). The body pose ( $\boldsymbol{\theta}_b$ ) is modeled using VPoser [71], a variational autoencoder that transforms the pose into a 32-dimensional latent vector  $Z$ . An L2 prior is enforced in this latent space:  $E_{\theta_b}(\boldsymbol{\theta}_b) = \|Z\|^2$ . For the hands, SMPL-X uses a low-dimensional PCA pose space such that  $\boldsymbol{\theta}_h = \sum_{n=1}^{|m_h|} m_{h_n} \mathcal{M}$ , where  $\mathcal{M}$  are principal components capturing finger pose variations and  $m_{h_n}$  are the corresponding PCA coefficients. Thus,  $E_{m_h}(m_h)$  is an L2 prior on the coefficients  $m_h$ . Additional priors include:  $E_\alpha(\boldsymbol{\theta}_b) = \sum_{i \in \{\text{elbows, knees}\}} \exp(\boldsymbol{\theta}_i)$ , which penalizes extreme bending of elbow and knee joints;  $E_{\theta_f}$ ,  $E_\beta$ , and  $E_\epsilon$  as squared L2 priors for face pose, body shape, and facial expression, respectively; and the collision term  $E_C$  prevents self-interpenetration by pulling intersecting vertices to the mesh surface. The optimization runs in multiple stages to prevent small body parts from dominating early iterations, starting from a mean pose and progressively refining the estimate.

### 2.2.3 Challenges for Sign Language Capture

Even though HPE has made significant progress, capturing SL poses (which we defined as SLC) presents several challenges:

1. **Handshape complexity:** hands have many degrees of freedom [5], making it difficult to capture the precise handshapes that are crucial in SL. However, an accurate reconstruction is critical for distinguishing minimal pairs in which only the handshape differs (*e.g.*, SCHOOL versus IMPOSSIBLE).
2. **Motion blur:** SL movements often involve rapid transitions that cause motion blur in recordings [89]. Heavy motion blur challenges even human observers to understand the handshape from the visual evidence of a frame.
3. **Self-touch:** SLs heavily rely on self-touch, particularly in Type X and Type 3 signs. Self-contacts are inherently difficult to reconstruct [66].
4. **Depth ambiguity and self-occlusion:** most HPE methods handle 2D image-plane estimation reasonably well but struggle with depth (*i.e.*, along the camera’s view-

ing direction). This limitation poses a major challenge for SLC, where self-contact often occurs along this axis: when a hand moves near the face or body, a frontal camera view alone typically cannot determine whether contact actually occurs, especially in the presence of self-occlusion.

5. **Non-manual features:** non-manual markers, including facial expressions, are integral to SL grammar and meaning, yet standard HPE methods often overlook or poorly capture them [46].

These challenges show that effective SLC needs tailored approaches.

## 2.2.4 Current Methods for Sign Language Capture

Prior work has approached SLC from multiple angles. Kratimenos *et al.* [46] compare recognition accuracy on Greek Sign Language [86] using features from raw RGB images, OpenPose [9] 2D skeletons, and SMPL-X bodies. The SMPL-X bodies, obtained with SMPLify-X, demonstrate the best results, highlighting the value of 3D data. How2Sign [17] provides 3D skeleton reconstructions for three hours of data captured in a Panoptic Studio [41]. While this approach offers some insights into SL motion, the skeletal representation lacks the surface details necessary to capture crucial phonological features such as the fine-grained facial expressions [63] available in models like SMPL-X. SignPose [47] infers a textured avatar from single RGB images using synthetic SL animations for training but has limited application, as it requires all upper body keypoints to be detected, which is an unrealistic expectation in SL videos.

These results demonstrate that while 3D body models like SMPL-X show promise for SL, existing methods do not fully address the SL-specific challenges described above. Notably, although linguistic research has long established the phonological structure of SLs, a significant research gap exists in using these linguistic insights as “side information” for improving SLC. This gap between linguistic theory and computer vision practice presents an opportunity for interdisciplinary approaches that formalize linguistic knowledge as computational priors. Additionally, the presence of self-touch events has been largely overlooked, despite their phonological significance.

## 2.3 Sensing Self-Contact

A wide range of sensing modalities has been explored to detect self-touch, or more generally self-contact, each offering different trade-offs in accuracy, intrusiveness, and scalability:

- **Vision:** RGB, stereo, and depth (RGB-D) camera systems have been widely used to automatically identify the beginning and end of self-contact events [14, 22, 36, 53, 62, 67]. These systems rely on pose tracking and motion analysis to infer

self-contact but often struggle to distinguish between contact and close proximity. Their performance further degrades under occlusions and variable lighting. Manual annotation of videos is sometimes used as a fallback [50] but is labor-intensive, inconsistent, and impractical for large-scale deployment.

- **Motion and proximity:** sensors such as accelerometers and electromyography units have been used to infer gestures and interactions involving self-touch [64]. These systems are compact and practical but do not directly detect contact, often mistaking motion for self-contact. Similarly, short-range radar [32] can detect hand proximity but encounters similar ambiguity when differentiating contact from close proximity.
- **Acoustic and ultrasound:** bioacoustic systems and ultrasound propagation techniques are another indirect way to detect contact events [35, 49, 65]. However, signal strength attenuates with increasing distance between the contact site and sensor [35, 65], or sensors must be placed on the body part where contact occurs [49]. Consequently, detecting self-contact across multiple body regions would require extensive instrumentation.
- **Tactile:** flexible e-skin technologies can directly measure contact location and pressure [93]. These systems offer high-fidelity data but must be worn over all potential contact areas, which reduces wearability and comfort. They are more suitable for use in controlled settings or for providing reliable pseudo-GT (pGT) for visual data in contact situations involving limited body parts rather than full-body scenarios [62].
- **Electrical:** electrical conductivity sensing can detect [97] or localize [96] skin contact but only in limited spatial regions due to the required proximity of electrodes. In contrast, electrical bioimpedance sensing offers a more scalable alternative; the Touché system [80] demonstrates bioimpedance-based classification of static poses, including within-hand and hand-to-head contact.

### 2.3.1 Electrical Bioimpedance Sensing

Besides its usage for classifying poses as in Touché [80], electrical bioimpedance sensing represents a promising approach for detecting the beginning and end of self-contact. Bioimpedance measures how strongly a biological medium opposes the flow of an alternating current. As such, the bioimpedance ( $Z$ ) between two points is defined as the ratio of voltage to current ( $Z = V/I$ ) [43]. As a complex number, bioimpedance combines resistance ( $R$ , real part) and reactance ( $X$ , imaginary part), expressed through two frequency-dependent components: magnitude ( $|Z| = \sqrt{R^2 + X^2}$ ) and phase angle ( $\angle Z = \tan^{-1}(X/R)$ ). The resistance component is primarily caused by total body water, which is moderately conductive, while the reactance component is mainly due to

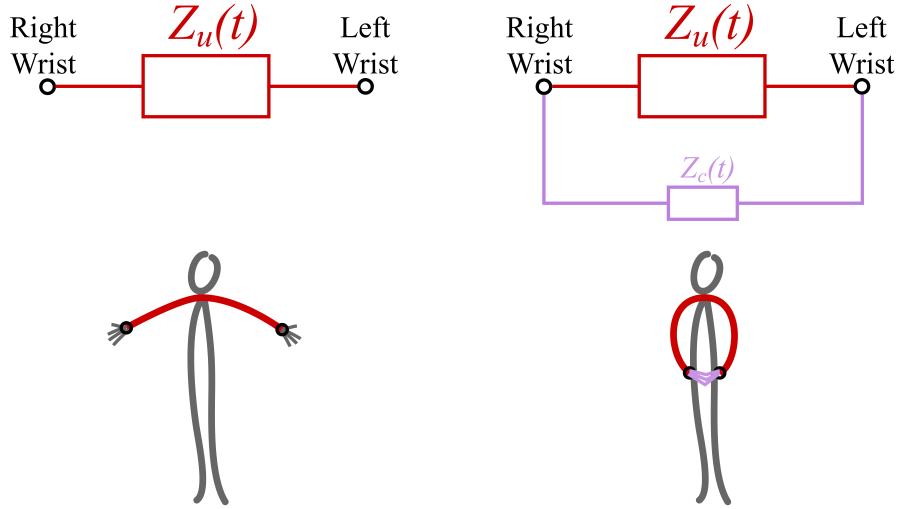


Figure 2.2: The bioimpedance of the user’s body can be measured from wrist to wrist ( $Z_u$ ). Forming a new electrical pathway between the wrists, such as when the hands contact each other, can be modeled by an additional bioimpedance component in parallel with the original one ( $Z_c$ ), as shown on the right. We define both impedance values as functions of time ( $Z(t)$ ) because the bioimpedance of a person changes continuously.

cell membranes acting as capacitors [68]. For practical implementation, small alternating currents (typically  $100\mu\text{A}$  to  $1\text{mA}$ ) at frequencies ranging from  $1\text{kHz}$  to  $100\text{kHz}$  are applied to the body. Higher frequencies are often preferred as they penetrate cell membranes more effectively [15].

The bioimpedance measured between two anatomical locations (such as the wrists) varies significantly between individuals due to differences in body composition and shape [90]. However, even within the same individual, it fluctuates due to factors including skin temperature, core temperature, body position, muscle contraction, exercise, hydration, and fasting state [3, 13, 20, 31, 60]. Most critically for self-contact detection, when two distinct body parts come into contact, a new electrical pathway is created through which current can flow in parallel with the default pathway through the body, as illustrated in Fig. 2.2. This new pathway leads to a measurable decrease in the total bioimpedance, following the formula for parallel impedances:

$$Z_{\text{total}} = \left( \frac{1}{Z_u} + \frac{1}{Z_c} \right)^{-1}, \quad (2.9)$$

where  $Z_u$  is the bioimpedance of the upper body measured between the wrists, and  $Z_c$  is the bioimpedance of the path formed by the new contact. For example, when a signer articulates THANKS, they begin the sign by briefly touching their lips. This momen-

tary contact creates a parallel electrical pathway, which temporarily lowers the measured bioimpedance.

Bioimpedance offers several key advantages that make it particularly suitable for detecting self-touch in SL:

- **Contact specificity:** the measurable changes in bioimpedance occur only when actual physical contact is made between body parts, providing a clear distinction between Type 0 (no contact) and Type X signs (with contact).
- **Temporal characterization:** different types of contact produce characteristic temporal patterns in the bioimpedance signal, enabling bioimpedance-based systems to distinguish between quick taps (*e.g.*, THANKS), repeated contacts (*e.g.*, NAME), and sustained holds (*e.g.*, PLEASE).
- **Immunity to visual occlusion:** bioimpedance sensing directly measures the electrical properties of the body, making it immune to the visibility issues that commonly affect visual systems.
- **Wearability and power efficiency:** advances in bioimpedance systems have led to highly miniaturized and energy-efficient designs [15]. These compact devices can now be integrated into ergonomic form factors capable of operating for extended periods on a single battery charge, making them suitable for continuous long-term use.

However, while bioimpedance sensing can be used to detect the time of a contact, it cannot precisely localize it without sensorizing large body areas.

## 2.4 Integrating Linguistics and Bioimpedance Sensing into Human Pose Estimation

The previous sections have established three key insights that motivate this dissertation's interdisciplinary approach. First, SL linguistics provides systematic knowledge about how signs are structured, particularly through Battison's classification and the critical role of self-touch in distinguishing sign types [4]. Second, current HPE methods, while advancing rapidly, face fundamental limitations when applied to SL due to handshape complexity and self-contact ambiguities. Third, bioimpedance sensing offers a promising solution for detecting the temporal characteristics of self-contact events. Importantly, these insights address complementary aspects of the SLC challenge: linguistic knowledge can constrain the hand pose estimation problem, vision-based HPE methods can provide spatial localization of the body parts in contact, and bioimpedance sensing can resolve the temporal aspect of the contact. Together, they address limitations that each modality alone cannot overcome. This dissertation explores this integration through a

progressive approach. Chapter 3 investigates how linguistic constraints can be formalized as computational priors and integrated into optimization-based HPE frameworks. Chapter 4 systematically characterizes bioimpedance sensing for self-touch detection, determining optimal parameters across individuals. Chapter 5 demonstrates how temporal contact information from bioimpedance can be integrated with spatial pose estimates from vision-based HPE methods during self-touch.

This interdisciplinary framework not only advances SLC capabilities but also establishes methodological principles for addressing similar challenges in other HPE applications where domain knowledge and complementary sensing modalities can overcome limitations of HPE methods alone.



# Chapter 3

## Integrating Linguistic Priors into Human Pose Estimation

**Note:** This chapter is based on Forte, Kulits, Huang, Choutas, Tzionas, Kuchenbecker, and Black’s article “Reconstructing Signing Avatars From Video Using Linguistic Priors,” which was published in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 2023, pp. 12791-12801* [24]. Some paragraphs in this dissertation’s Abstract, Introduction, Background, and Discussion are also adapted from this publication.

As discussed in the previous chapter, HPE methods struggle with SL for several reasons, including the complexity of hand movements and the high number of degrees of freedom present in hands [5]. Additionally, motion blur induced by fast motions [89], especially in low-resolution videos, causes hand pose to be unrecognizable in many frames (see Fig. 3.1). This chapter focuses on improving hand pose estimation by developing linguistic-based priors that help disambiguate hand poses in SL videos; this is a novel use of linguistic “side information” to improve HPE. Note that in HPE, what we refer to as “hand pose” corresponds to the “handshape” parameter in SL linguistics; both terms describe the configuration of the fingers.

Based on hand poses and movements, Battison [4] defines five main linguistic classes that represent all signs (refer to Sec. 2.1.2 for more information). We build on his work to define eight classes; we combine all one-handed signs into class 0, while two-handed signs are arranged in classes 1, 2, or 3, depending on how the non-dominant hand participates in the articulation of the sign. We further divide each of these four classes into two subclasses depending on whether the pose of the active hand(s) changes during the articulation of the sign. Our subdivision captures hand-pose dynamics, reflecting distinctions that are important for computer vision but not explicitly emphasized in linguistic analysis.

We introduce two class-dependent SL linguistic constraints that systematically describe our eight classes: hand-pose symmetry and hand-pose invariance. Under Battison’s SL symmetry condition [4], when both hands actively move, the articulation of the

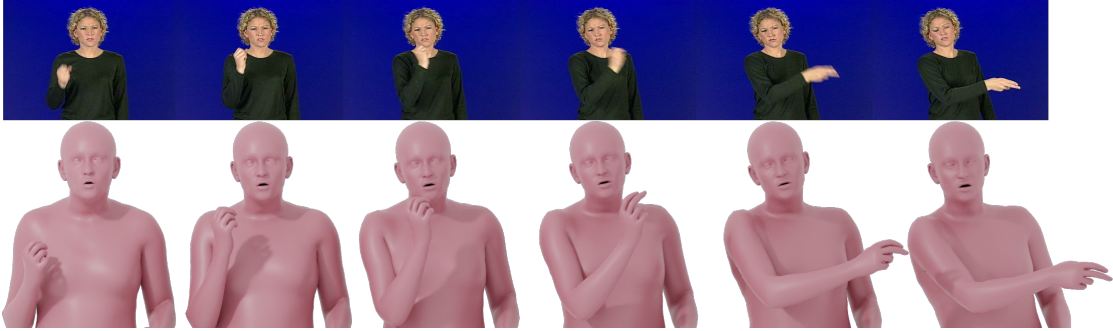


Figure 3.1: Given a monocular, in-the-wild video of a SL sign, SGNify automatically reconstructs a 3D body with accurate hand pose, facial motion, and body pose. Note that motion blur obscures the finger articulations in several video frames; this is a common problem. Our novel linguistic priors help enable 3D reconstruction despite such image degradation.

fingers must be identical. We formalize this concept as a regularization term that encourages the pose of the two hands to be similar for such signs. Coupling the hand poses in this way effectively increases the image evidence for a pose, which improves estimates for challenging videos. Our invariance constraint uses the observation that hand pose is either static or transitions smoothly from one pose to another during the articulation of the sign; other significant changes to hand pose are not common in SL. Specifically, we extract a characteristic “reference pose sequence” (RPS) to describe each local hand pose during the sign articulation, and we penalize differences between the RPS and the estimated hand pose in each frame. These two priors of hand-pose symmetry and hand-pose invariance are universally applicable to all SLs.

Our novel hand-pose constraints are formulated to be incorporated into the objective function of optimization-based methods or into the loss function for training a neural network regressor. Due to the limited training data, we use an optimization-based approach and build on the SMPLify-X method [71] (refer to Sec. 2.2.2 for more information). To appropriately incorporate our terms into the objective function, we need to know the class of the sign. We train a simple model that extracts features from the raw video and determines the class to which the depicted sign belongs. While SMPLify-X is a good foundation for the hands and body, we find that it does not capture expressive facial motions well. Consequently, we use a more expressive face regressor, SPECTRE [23], to capture the face parameters. We call our method SGNify.

To quantitatively evaluate SGNify, we capture a native DGS signer with a frontal RGB camera synchronized with a 54-camera Vicon motion capture system and recover ground-truth meshes from the Vicon markers [56]. We run SGNify on the RGB video and compute 3D vertex-to-vertex (V2V) error between our resulting avatars and the ground-truth meshes. We find that SGNify reconstructs SMPL-X meshes more accurately than existing methods.

We conduct a perceptual evaluation in which we present proficient signers with a video of either an estimated SMPL-X avatar or the real-person source video and task them with identifying the sign being performed. Participants also rate their ease in recognizing the sign and the naturalness of the articulation. Our results show that SGNify reconstructs 3D signs that are as recognizable as the original videos and consistently more recognizable, easier to understand, and more natural than the existing HPE methods. We also evaluate SGNify in a multi-view setting and on continuous signing videos. Despite not being designed for the latter, SGNify captures the meaning in continuous SL.

SGNify represents a step toward the automatic reconstruction of natural 3D avatars from SL videos. Our key contribution is the introduction of novel linguistic priors that are universal and helpful to constrain the problem of hand pose estimation from SL video. SGNify is designed to work on video from different SL dictionaries across languages, backgrounds, people, trimming, image resolution, and framing. This capability is critical to capture 3D signing at scale, which will enable progress on learning SL avatars. Our code and motion-capture dataset are available for research purposes at <https://sgnify.is.tue.mpg.de>

## 3.1 SGNify

SGNify is an offline method for reconstructing 3D body shape and pose of isolated SL signs from monocular RGB video. SGNify centers around a key insight: SL signs follow universal linguistic rules that can be formulated as class-specific priors and used to improve hand pose estimation. Our full pipeline is shown in Fig. 3.2.

### 3.1.1 SMPLify-SL: Baseline for Sign Language Video

Our baseline method builds on SMPLify-X [71], which estimates SMPL-X [71] parameters from RGB images.

To create a strong baseline, we extend SMPLify-X to video by adapting it in the following ways: (1) We cope with the upper-body framing typical of SL videos by changing the heuristic used for camera initialization and the estimation of the out-of-view lower-body joints. (2) Since human motion is locally smooth in time, we initialize  $\theta_t \in \mathbb{R}^{|\theta|}$  with  $\theta_{t-1}$  and include a zero-velocity loss on the hands and body to encourage smooth reconstructions. (3) We estimate shape parameters ( $\beta$ ) over multiple frames by taking the median of the parameter estimates and not-optimizing them during the per-frame reconstruction. (4) To better capture the frequent hand-hand and hand-body interactions (mainly with the face and the chest), we employ the more robust self-contact loss of Müller *et al.* [67] instead of the original SMPLify-X interpenetration term. (5) For each frame, we pre-compute the facial expressions ( $\psi$ ) and jaw poses with SPECTRE [23]. These parameters are substituted into SMPL-X at the end of the optimization. SPECTRE

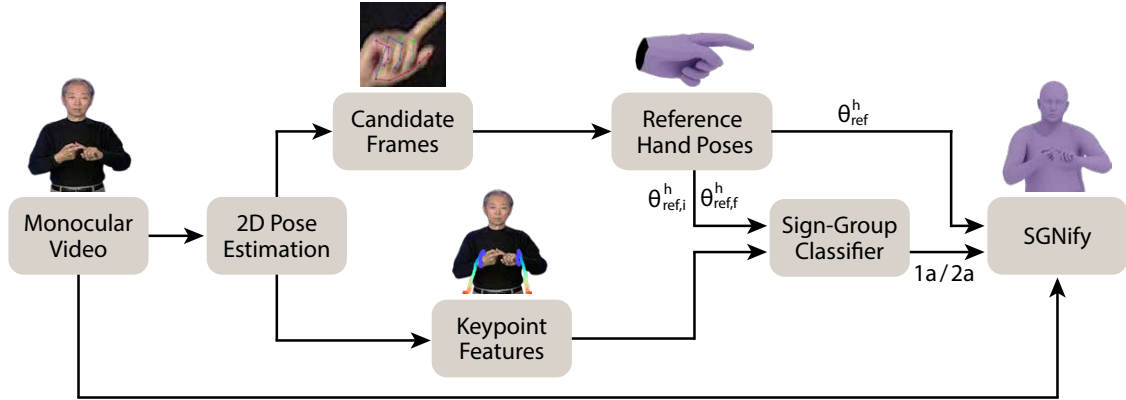


Figure 3.2: Given a video of a SL sign as input, our method preprocesses the data to first extract 2D keypoints. The hand keypoints are used to select candidate frames for estimating the reference hand poses ( $\theta_{ref,i}^h$ ,  $\theta_{ref,f}^h$ , and, for static hand poses, also  $\theta_{ref}^h$ ). The initial and final reference hand poses ( $\theta_{ref,i}^h$  and  $\theta_{ref,f}^h$ ), together with wrist-keypoint features detected across the sequence, are then fed into our sign-group classifier, which classifies signs into groups based on linguistic rules universally applicable to SL [4]. Using the predicted group labels and the relevant reference hand poses, SGNify applies the appropriate linguistic constraints to improve SL 3D hand-pose estimation, especially when the video frame is ambiguous.

can be swapped for any method whose expression parameters are consistent with those of SMPL-X, *e.g.*, EMOCA [11]. We denote the baseline SMPLify-SL.

### 3.1.2 Linguistic Constraints

Optimization- and regression-based human pose estimation methods struggle on SL video, particularly with the estimation of hand pose. We address this challenge by formulating linguistic constraints as additional losses on hand pose and integrating them into the SMPLify-SL objective function. First, we adapt the five sign-classification and morpheme-structure conditions introduced for American Sign Language (ASL) by Battison [4] to divide signs into four primary classes:

- **Class 0** one-handed signs in which only the dominant hand articulates the sign.
- **Class 1** two-handed signs in which both hands are active. They share the same poses and perform the same movement in a synchronous or alternating pattern. This class includes all signs that follow Battison’s symmetry condition [4].
- **Class 2** two-handed signs in which the dominant hand is active, the non-dominant hand is passive (its position and pose do not change during the articulation of the sign), and the two hands have the same initial pose.

- **Class 3** two-handed signs in which the dominant hand is active, the non-dominant hand is passive, and the two hands have different hand poses. All signs in this class follow Battison’s dominance condition [4].

We further divide each class into two subclasses: *subclass a* contains signs in which the hand pose of the active hand(s) does not change throughout the articulation of the sign (“static”), and *subclass b* contains all signs in which the hand pose changes (“transitioning”).

Note that the division into these classes is not limited to ASL; Eccarius *et al.* [18] show that the phonological and prosodic properties of ASL can be successfully transferred to other sign-language lexicons.

We then convert these linguistic classes into two 3D pose constraints: hand-pose symmetry and hand-pose invariance. Signs in the same class share the same constraints (see Tab. 3.1).

### Hand-Pose Symmetry

We encourage the left and right hand poses to match for the relevant classes (classes 1a, 1b, and 2a in Tab. 3.1):

$$L_s = \lambda_s \|\theta_t^r - r(\theta_t^l)\|_2^2, \quad (3.1)$$

where  $\theta_t^r$  is the finger articulation of the right hand, and  $r(\theta_t^l)$  is a reflection function to represent the articulation of the fingers of the left hand as if it were a right hand. This loss penalizes differences in finger poses between the hands.

### Hand-Pose Invariance

Each sign has a characteristic reference hand pose sequence (RPS). The RPS defines the hand pose that we expect at each time  $t$  during the articulation of the sign. The hand-pose-invariance constraint penalizes differences between the reference hand pose  $\theta_{ref,t}^h \in \text{RPS}^h$  and the estimated hand pose  $\theta_t^h$ :

$$L_i^h = \lambda_i \|\theta_{ref,t}^h - \theta_t^h\|_2^2, \quad (3.2)$$

where  $h$  represents either the left or the right hand.

Throughout each sign, the hand pose either stays static or transitions between two poses. When static, only one hand pose,  $\theta_{ref}^h$ , is representative of the RPS. Signs where the hand pose is transitioning are characterized by two reference hand poses,  $\theta_{ref,i}^h$  and  $\theta_{ref,f}^h$ , corresponding respectively to the initial and final poses. We interpolate  $\theta_{ref,i}^h$  and  $\theta_{ref,f}^h$  with spherical linear interpolation [82] to obtain intermediate poses. We presently do not consider signs with repeated hand-pose transitions, *e.g.*, STORY in ASL, which occur in a small percentage of signs ( $\sim 3\%$ ).






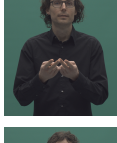
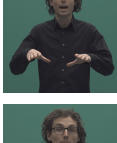
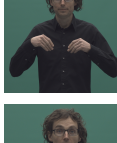
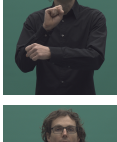
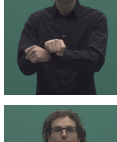
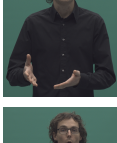
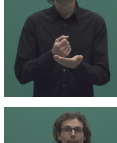

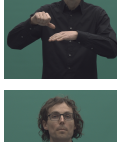

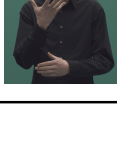
Initial Hand Pose	Final Hand Pose	Class	Hand-Pose Symmetry	Hand-Pose Invariance Dominant	Hand-Pose Invariance Non-dominant
		0a	✗	static	✗
		0b	✗	transitioning	✗
		1a	✓	static	static
		1b	✓	transitioning	transitioning
		2a	✓	static	static
		2b	✗	transitioning	static
		3a	✗	static	static
		3b	✗	transitioning	static

Table 3.1: Linguistic constraints defining the eight sign classes. The images are representative of our eight sign classes. The videos of these signs appear in the supplemental video at <https://sgnify.is.tue.mpg.de>.

The full objective function of SGNify is thus:

$$E(\beta, \theta, \psi) = E_J + \lambda_{\theta_b} E_{\theta_b} + \lambda_{m_h} E_{m_h} + \lambda_{\alpha} E_{\alpha} + E_O + \lambda_P E_P + \lambda_A E_A + L_s + \sum_{h \in \{r, l\}} L_i^h + \lambda_t L_t + \lambda_{st} L_{st}, \quad (3.3)$$

For more details on the first four terms, please refer to Sec. 2.2.2 or the original paper of SMPLify-X [71].  $E_O$  is a bone-orientation term, which measures how well the directions of bones (*i.e.*, the vectors between parent and child joints) align between the target points (detected landmarks) and the estimated SMPL-X joints, ensuring realistic bone orientations. For more details about this term, please refer to the original paper of RICH [40].  $E_P$  and  $E_A$  are used to prevent self-interpenetration. When self-contact occurs, the  $E_P$  term pushes vertices that are inside the mesh to the surface, and  $E_A$  aligns the surface normals of the vertices in contact. For more details, please refer to the original paper of TUCH [67].

We add  $L_s$  and  $L_i^h$  to enforce our linguistic constraints:  $L_s$  represents the symmetry constraints, and  $L_i^h$  the hand-pose invariance of the right ( $r$ ) and left ( $l$ ) hands. We also add a temporal loss  $L_t$  on the body- and hand-pose vectors and a standing loss  $L_{st}$  to penalize deviations from a standing pose when none of the feet keypoints are detected; specifically, this penalization is applied to the joints below the pelvis and to the spine.

Finally, each  $\lambda$  denotes the influence weight of each loss term. For more details on the exact  $\lambda$  values and insights on the full SGNify objective, please see the code, which can be reached from the project URL.

We optimize our objective function using the trust-region Newton conjugate gradient method [69]. Note that we do not optimize the shape  $\beta$  or the facial expressions  $\psi$ .

### 3.1.3 Automatization

To work fully automatically, SGNify must 1) estimate the poses needed to enforce the hand-pose-invariance constraint and 2) classify which sign group is present in a video sequence (see Fig. 3.2).

To estimate the reference hand poses ( $\theta_{ref}^h$ ,  $\theta_{ref,i}^h$ , and  $\theta_{ref,f}^h$ ), our method selects candidate frames in the core part of the sign using hand-keypoint detection confidences, and it uses SMPLify-X (adapted to SL cropping) to reconstruct a preliminary 3D hand pose for each candidate frame. With static hand poses,  $\theta_{ref}^h$  is obtained by taking the average hand poses of these candidates. With transitioning hand poses, the core part of a sign is divided into two intervals, and  $\theta_{ref,i}^h$  and  $\theta_{ref,f}^h$  correspond to the average hand poses of the candidate frames in the first and second intervals, respectively. When articulating an isolated sign, signers start and end in a rest pose. SGNify identifies the beginning and end of the sequence based on when the hands begin to move. After automatic trimming, the initial and final frames of the sequence show the transition from the rest pose to the pose(s) characteristic of the sign. We observe that the transition from the rest pose to

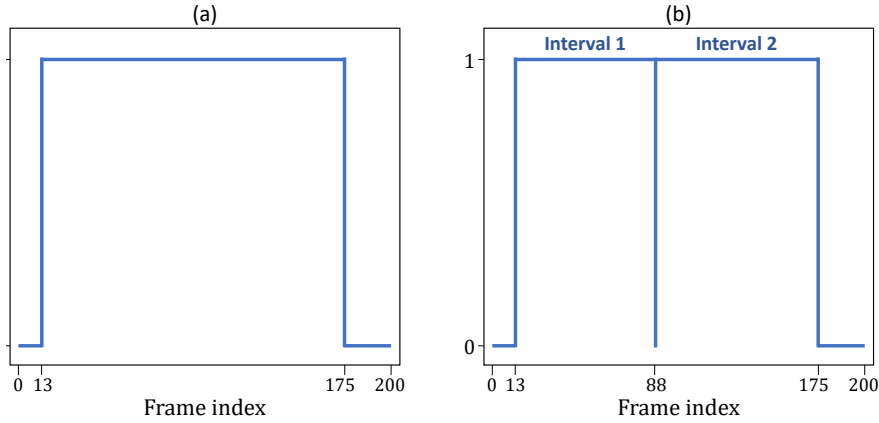


Figure 3.3: We consider an example sequence of 200 frames. (a) Static hand: Frames whose value on the y-axis is 1 are candidates for identifying  $\theta_{ref}^h$ . (b) Transitioning hand and input features for the sign-group classifier: The first interval shows candidates for  $\theta_{ref,i}^h$ , and the second one for  $\theta_{ref,f}^h$ .

the core part of the sign usually happens around  $t = 0.5 * T/8$ , and the transition from the sign to the rest pose typically occurs around  $t = 7 * T/8$ , where  $T$  is the number of frames in the motion sequence. As a result, we assume the core part of a sign to happen between  $0.5 * T/8 < t < 7 * T/8$ . Figure 3.3a shows the frames during which the transitions from/to the rest pose happen (indicated with 0) and the frames during which the sign is articulated (indicated with 1) for a sample trimmed recording containing 200 frames. To identify the two key poses representing the initial and final hand poses ( $\theta_{ref,i}^r$  and  $\theta_{ref,f}^r$ ), we consider two different intervals; we expect to see the first hand pose at the beginning of the sequence (first interval shown in Fig. 3.3b) and the second hand pose at the end (second interval shown in Fig. 3.3b).

The constraints applied to each sign depend on its sign group; we have six sign groups because classes 1a & 2a share the same constraints, as do 2b & 3b (Tab. 3.1). There is insufficient paired data to train a CNN classifier, so we use an intuitive and interpretable decision tree trained on extracted 2D and 3D pose features. Our features are invariant to the handedness of the signer and include:

- 1) the minimum of the normalized maximum height differences of each wrist across the sequence (normalization by the median nose–neck distance, as a proxy for body scale):  $\min(\{w_r\}_{\max} - \{w_r\}_{\min}, \{w_l\}_{\max} - \{w_l\}_{\min})$ , where  $w_l$  and  $w_r$  are the heights of the left and right wrists, respectively.
- 2) the cosine distance between the initial poses of each hand:  $\text{CosDist}(\theta_{ref,i}^r, \theta_{ref,i}^l)$ ,
- 3) the maximum of the two cosine distances between each initial and final hand pose:  $\max(\text{CosDist}(\theta_{ref,i}^r, \theta_{ref,f}^r), \text{CosDist}(\theta_{ref,i}^l, \theta_{ref,f}^l))$ .



Figure 3.4: Samples frames reconstructed by SGNify. The input videos are from diverse multilingual datasets. Row 1: GDAŃSK in PJM (class 1a) from the dataset [52] used for training our sign-group classifier. Row 2: BLUME in DGS (class 3b) from our captured dataset used in the quantitative evaluation. Row 3: DOLL in ASL (class 0a) from the dataset [85] used in the perceptual study. Row 4: ARCHERY in South African SL (SASL, class 2a) from the in-the-wild dataset [75].

We train our sign-group classifier on the over 3,000 videos from the Corpus-based Dictionary of Polish Sign Language (CDPSL) [52]; these are annotated with HamNoSys [33], an alphabetic system describing signs on a mostly phonetic level. Row 1 of Fig. 3.4 shows a sample frame from CDPSL. This dataset is not used in our quantitative analysis or perceptual study. Using our EBNF grammar (see Sec. 2.1.3), we parse HamNoSys [33] annotations on CDPSL [52] and extract labels to train our sign group classifier. We assign our classes to the clips as follows:

- **Class 0a** There is one *handshape\_block* nonterminal and no SYMMETRY terminal is present.
- **Class 0b** There are two *handshape\_block* nonterminals, the two *handshape\_block* nonterminals are not equal, a HAMREPLACE terminal is present, and no SYMMETRY or REPEAT terminals are present.
- **Class 1a** There is one *handshape\_block* nonterminal and a SYMMETRY terminal is present.
- **Class 1b** There are two *handshape\_block* nonterminals, they are not equal, a HAMREPLACE terminal is present, and no SYMMETRY or REPEAT terminals are present.
- **Class 2a** There are two *handshape\_block* nonterminals, they are equal, they fall within a *par* nonterminal, and no SYMMETRY terminal is present.
- **Class 2b** There are three *handshape\_block* nonterminals, the first two are equal, a HAMREPLACE terminal is present, and no SYMMETRY or REPEAT terminals are present.
- **Class 3a** There are two *handshape\_block* nonterminals, they are not equal, they fall within a *par* nonterminal, and no SYMMETRY terminal is present.
- **Class 3b** There are three *handshape\_block* nonterminals, the first is not equal to the second, a HAMREPLACE terminal is present, and no SYMMETRY or REPEAT terminals are present.

Note that the SYMMETRY parameter from HamNoSys refers to Battison’s symmetry condition [4], which also includes the signer’s arm movement and not only the hand pose; in contrast, our symmetry constraint applies only to hand pose.

## 3.2 Dataset

To quantitatively evaluate SGNify as a viable method for SLC, we collected motion-capture data with ground-truth SMPL-X bodies articulating signs. Our dataset represents

<b>Class</b>	0a	0b	1a	1b	2a	2b	3a	3b
<b># Signs</b>	12	3	14	3	11	2	10	2

Table 3.2: Number of signs captured for each class.

the first publicly available expressive full-4D capture of isolated SL signs. The experimental procedure was reviewed by the ethics council of the University of Tübingen without objections or remarks (709/2021B02).

In consultation with a Deaf DGS teacher and a DGS interpreter, we defined a DGS corpus consisting of 57 isolated signs. The selected signs cover a wide range of challenges for SLC, such as self-contact and self-occlusion. Table 3.2 summarizes the number of signs collected for each of the eight classes. Signs of *subclass b* are less common, and this is reflected in our corpus.

We captured a native right-handed DGS signer with a Vicon mocap system at 120 fps, synchronized with a frontal  $4112 \times 3008$  RGB camera at 60 fps, framing an upper-body view as typically found in SL video. The hands start and end at rest at the signer’s sides, and each sign lasts between 1.7 and 3.5 seconds after trimming. In total, our dataset comprises 16,608 mocap frames and 8,304 RGB frames. To obtain ground-truth SMPL-X meshes, we scanned the participant in a 4D body scanner in several poses. The SMPL-X mesh was registered to these scans and averaged to obtain a personalized body-shape mesh. MoSh++ was then used to fit this mesh to the mocap markers [56]. Marker-based mocap is useful for evaluating ground truth but is not practical for SLC at scale: it is expensive and requires expertise, the reflective markers attached to the signer can influence contact-heavy motions, and processing the resulting data is time-consuming. If our monocular method can approach the performance of mocap, it will be widely applicable.

## 3.3 Experiments

### 3.3.1 Quantitative Evaluation

We quantitatively evaluate SGNify, compare it with standard HPE methods, and quantify the improvement derived from each linguistic constraint. To emulate in-the-wild data, which might have very low resolution, low framerate, and an occluded lower body, we pre-processed our high-quality video data to a resolution of  $514 \times 300$  at 30 fps, and we cropped the input images above the pelvis (see Row 2 of Fig. 3.4). We used the synchronized meshes captured from the observed Vicon markers [56] as ground truth for evaluation. Since all tested methods estimate SMPL-X meshes with the same topology, we compute the mean per-vertex error (TR-V2V) by considering the vertices above the pelvis. The prefix “TR” means that we translationally align the mesh reconstructed for

each frame with the ground truth; *i.e.*, we center the meshes before computing these errors. Since the starting and ending transitions are not part of the sign itself, we manually annotate the expressive central portion of each sign from the raw videos and compute the quantitative results on only these central frames (in total, 2,872 RGB frames).

Table 3.3 shows the mean TR-V2V error across the 57 signs for four methods and three body regions. The column labeled “Upper Body” reports the error computed when considering the hands and the upper body (vertices above the pelvis). We include the head but not the face because our mocap system struggles to reconstruct face details using only 27 markers. We separately report the TR-V2V errors for the left and right hands because of the central role they play in SL. Figure 3.5 illustrates the subsets of vertices selected for the quantitative evaluation. This experiment compares SGNify with FrankMocap [78], PIXIE [21], PyMAF-X [94], and our baseline SMPLify-SL. SGNify achieves the lowest error for the upper body and both hands, compared with standard HPE methods.

Tables 3.4 and 3.5 show the improvements derived from each linguistic constraint. Table 3.4 reports the mean TR-V2V in symmetric signs (classes 1a, 1b, and 2a) when using no constraints, *i.e.*, SMPLify-SL (Row 1), only one constraint (Rows 2 and 3), or both linguistic constraints (Row 4). When one linguistic constraint is used, the TR-V2V of both hands decreases. We observe that the symmetry constraint has the overall greater effect of the two. When the non-dominant hand is passive, it is often rotated at an angle difficult to capture from a frontal camera; this behavior might explain the greater effect of the symmetry constraint on the left hand. When both constraints are applied, we observe a more substantial decrease in the TR-V2V error of the left hand and a slight increase in the error of the right hand compared to when only the hand-pose invariance is applied to the right hand. We believe this happens because the symmetry constraint enforces symmetry between the two hands without knowing which hand reconstruction is more correct, and, for the same reason as above, detecting an accurate RPS for the non-dominant hand is often more challenging; overall, however, using both constraints greatly benefits the reconstruction. In Tab. 3.5, we separately evaluate performance on signs that do not present symmetry between the two hands (classes 0a, 0b, 2b, 3a, and

Method	Upper Body	Left Hand	Right Hand
FrankMocap [78]	78.07	20.47	19.62
PIXIE [21]	60.11	25.02	22.42
PyMAF-X [94]	68.61	21.46	19.19
SMPLify-SL	56.07	22.23	18.83
SGNify	<b>55.63</b>	<b>19.22</b>	<b>17.50</b>

Table 3.3: Evaluation on our ground-truth mocap dataset: mean TR-V2V error (mm) for five methods and three body regions.

Method	Sym	Inv	Left Hand	Right Hand	Both Hands
SMPLify-SL	✗	✗	20.30	18.78	19.54
SGNify	✓	✗	18.44	18.39	18.41
SGNify	✗	✓	19.76	<b>17.16</b>	18.46
SGNify	✓	✓	<b>17.72</b>	17.29	<b>17.50</b>

Table 3.4: Evaluating how the linguistic constraints of symmetry and invariance affect mean TR-V2V error (mm) in symmetric signs.

Method	Inv	Left Hand	Right Hand	Both Hands
SMPLify-SL	✗	26.09	18.89	20.50
SGNify	✓	<b>22.22</b>	<b>17.70</b>	<b>18.60</b>

Table 3.5: Evaluating how the linguistic constraint of hand-pose invariance affects mean TR-V2V error (mm) in asymmetric signs.

3b). As expected, the invariance prior included in SGNify improves the performance on all metrics. These quantitative results indicate that our linguistic constraints improve the reconstructions.

### 3.3.2 Qualitative Results

As previously mentioned, SGNify is designed to work on video from different SL dictionaries across languages, backgrounds, people, trimming, image resolution, and framing. Figures 3.6 and 3.7 show some of these examples and their SGNify reconstructions. In particular, Fig. 3.6 contains examples from the Real SASL [75] and CDPSL [52] datasets, and Fig. 3.7 from The American Sign Language Handshape Dictionary [85] and our collected DGS dataset (see Sec. 3.2).

### 3.3.3 Perceptual Study

We conduct an online perceptual study to 1) compare SGNify with the best-performing standard HPE method for SL and 2) evaluate the improvement derived from the linguistic constraints (SGNify vs. SMPLify-SL). Even though FrankMocap has a lower hand error than PyMAF-X and PIXIE, it has significant errors in the lower body when the full body is not seen, and higher errors in the upper body; these errors greatly affect the perceptual experience, making it unsuitable for the SL task, as already noticed in [47]. We thus use PyMAF-X since its hand-pose estimates are more accurate than PIXIE (see Tab. 3.3).

Approval of all experimental procedures for this study was granted by the Ethics Council of the Max Planck Society under the Haptic Intelligence Department’s framework

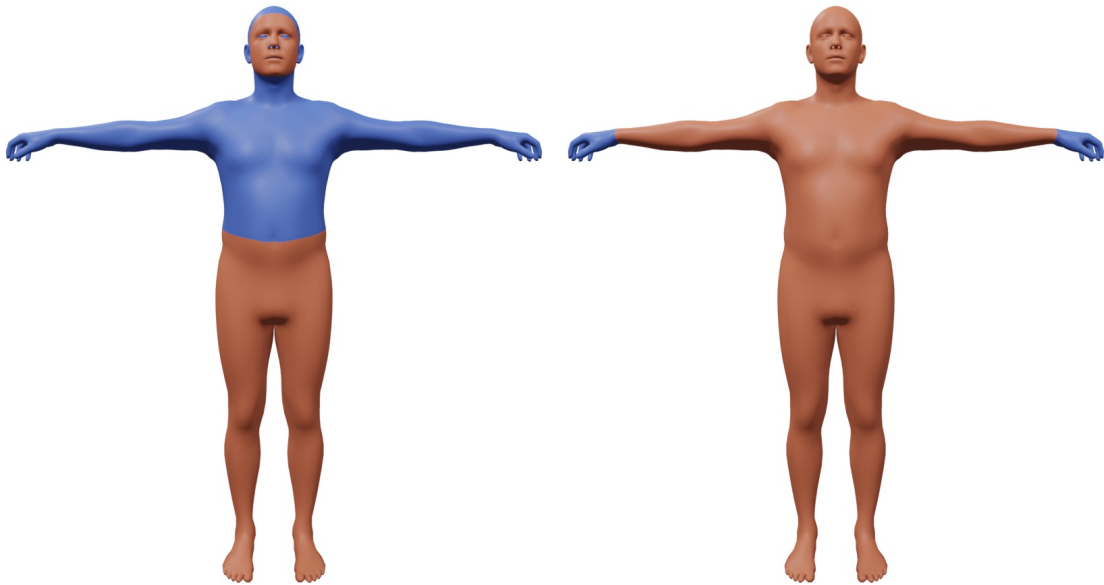


Figure 3.5: Blue vertices are used to calculate vertex error metrics, while red vertices are ignored. The left image shows the vertices used for the column of quantitative results labeled “Upper Body,” *i.e.*, upper-body vertices. The right image shows the vertex subsets for the left and right hands. Best viewed in color.

agreement under protocol number F027A. No participants are employed by our institution, and all are compensated 20 USD for their time.

Our study involves 20 adult participants who all stated that they have an advanced level of proficiency (expert level) in ASL. We discarded responses from other potential participants who did not correctly recognize at least 70% of the signs presented via real-person video. 15 of the final 20 participants (75%) are Deaf, and one participant is left-handed. Participant ages are  $46.75 \pm 11.78$ .

We used SGNify, SMPLify-SL, and PyMAF-X to reconstruct avatars from 50 videos taken from The American Sign Language Handshape Dictionary [85]; see Row 3 of Fig. 3.4 for an example. After responding to demographic questions, each participant evaluates the same six training videos to calibrate their responses to the quality of the presented reconstructions. We divide the remaining 44 signs into four batches (real-person RGB video, SGNify, SMPLify-SL, and PyMAF-X) that are balanced by sign class and sign frequency. We include real-person video to obtain an upper-bound on recognition performance and, as mentioned, to filter participants. We assign each participant to one of four surveys. Each survey contains all 44 test signs, and the four methods are rotated through the four sign batches across surveys. We further shuffle the questions in each survey for each user.

The participant enters the sign they believe the avatar or the real person is articulating in each video. They rate their ease in recognizing the presented sign using a visual analog

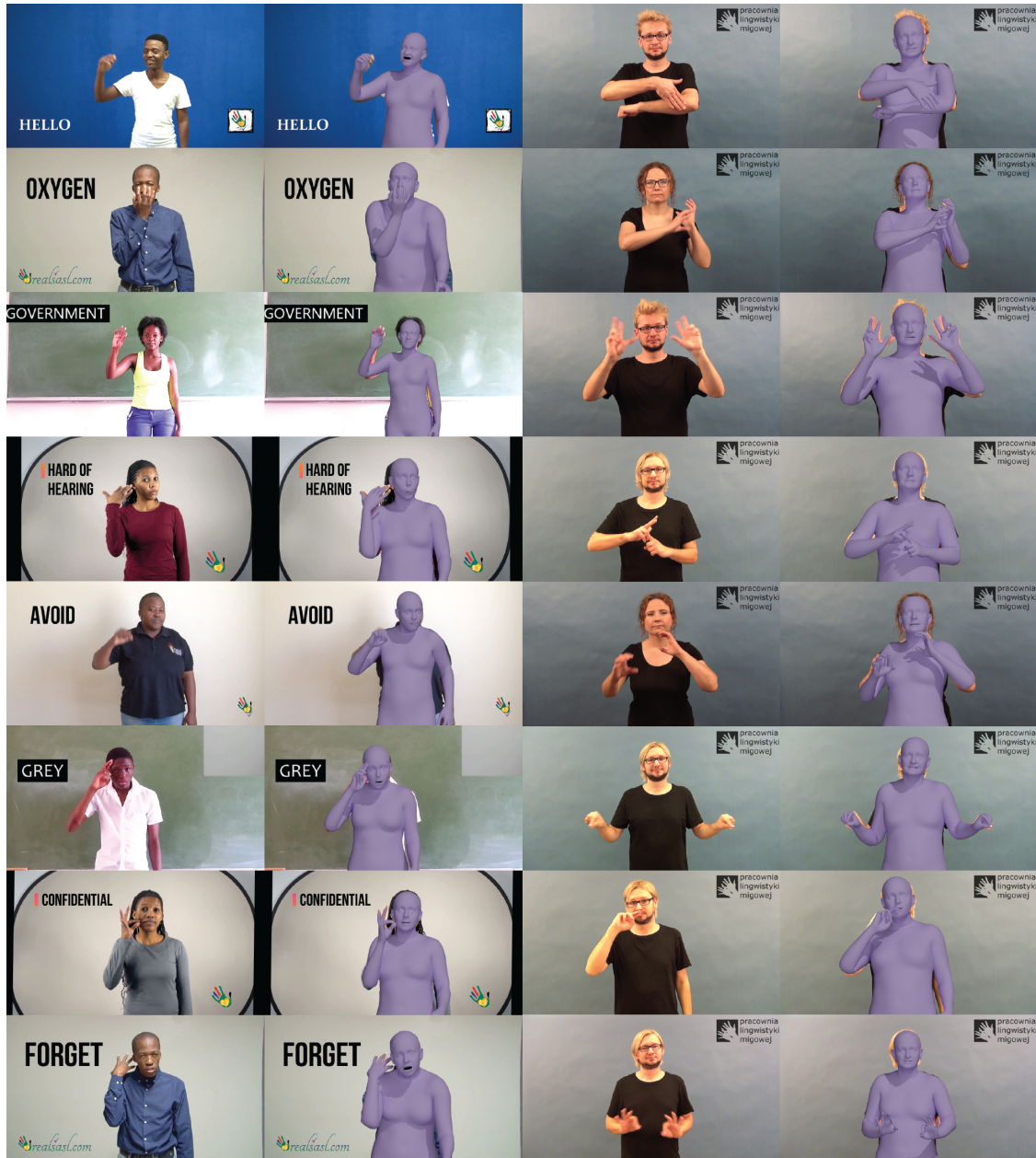


Figure 3.6: Additional examples on the Real SASL and CDPSL sign-language dictionaries, showing sample frames reconstructed by SGNify.



Figure 3.7: Additional examples on The American Sign Language Handshape Dictionary and our captured dataset, showing sample frames reconstructed by SGNify.

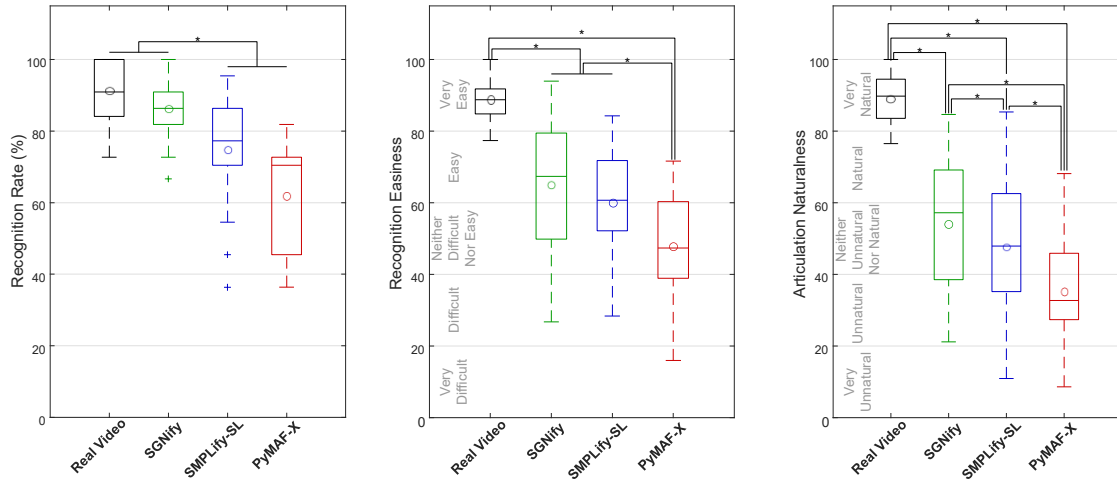


Figure 3.8: Box plots of the data from the perceptual study. The circle shows the mean, and the central line shows the median. The top and bottom box edges show the interquartile range (IQR), and the whiskers encompass the range up to 1.5 times the IQR. Outliers are marked with +. Statistically significant pairwise differences are indicated with a line and a  $\star$ . Left: Rate at which participants recognized signs presented with each of the four methods. *Real Video* and *SGNify* achieve significantly higher recognition rates than both *SMPLify-SL* and *PyMAF-X*. Center: Average easiness ratings participants assigned to recognizing signs presented by the four methods. All pairwise combinations except *SGNify-SMPLify-SL* are significantly different. Right: Average articulation naturalness ratings participants assigned to the four methods. All pairwise combinations are significantly different.

scale (VAS) ranging from 0 to 100 with five standard labels ranging from “very difficult” to “very easy.” They also evaluate the naturalness of the articulation on a VAS-labeled scale from “very unnatural” to “very natural.” Participants are able to replay each video. We provide additional space for comments for each sign and at the end of the study. The self-reported participant carefulness is  $84.45 \pm 11.44$  on a scale from 0 to 100.

The sign annotations submitted by participants are graded as either “incorrect” (no credit), “partially correct” (half credit), or “completely correct” (full credit). We calculate the rates at which each participant recognized the signs for each method. We visualize the resulting  $20 \times 4$  matrix of recognition rates with box plots in the left plot of Fig. 3.8. Participants recognize signs in real-person videos with an average accuracy of 90.9% and signs reconstructed by SGNify with 86.2% accuracy. Signs reconstructed with SMPLify-SL and PyMAF-X are recognized less accurately, at 74.8% and 62.0%, respectively. We evaluate the statistical significance of these recognition rates. Some distributions failed a Shapiro-Wilk normality test, so we used the non-parametric Friedman test which shows that the method used (real video, SGNify, SMPLify-SL, or PyMAF-X) has a statistically significant effect ( $p < 0.001$ ) on recognition rate. Pairwise compari-



Figure 3.9: Sample frames from the four methods presented in the second perceptual study: real video, the solid purple avatar from the first study, the same avatar wearing a black long-sleeved t-shirt, and a fully textured human character.

son with Wilcoxon signed-rank tests and a Bonferroni post-hoc correction reveal that the average sign recognition rate with real video and SGNify are both significantly higher than SMPLify-SL and PyMAF-X. Importantly, sign recognition rates with real video and SGNify are not significantly different from one another.

The central plot of Fig. 3.8 shows the participants’ perceived easiness in recognizing the sign. These four distributions passed the normality test and were analyzed with a one-way repeated-measures ANOVA with a Bonferroni post-hoc correction for pairwise comparisons. As expected, real videos are perceived to be significantly easier to recognize than the three reconstruction methods. However, signs reconstructed with SGNify and SMPLify-SL are significantly easier to recognize than those by PyMAF-X.

The right plot of Fig. 3.8 shows the participants’ perceived naturalness of the articulation of the signs. All four distributions passed the normality test, so they were analyzed in the same way as perceived easiness. All methods are statistically different one from another, with real video receiving the highest naturalness ratings, followed by SGNify, SMPLify-SL, and PyMAF-X. Nearly all participants reported in the comments that they want avatars with more expressive face motions and smoother hand and body motions.

These results show that SGNify outperforms the best-performing standard HPE method, PyMAF-X, in all qualitative metrics. Incorporating SL linguistic priors into our baseline SMPLify-SL yields statistically significant perceptual improvements, as judged by 20 expert signers. Most importantly, SGNify achieves a high sign-recognition rate that is not different from that of real-person video on the selected set of ASL signs.

A few participants suggested having clothed (rather than “nude”) avatars and introducing more human-like features. We thus conduct a second perceptual study to see whether these recommendations have a positive impact on the intelligibility of SGNify’s reconstructions. Thirteen participants from the first study participated in this follow-up; the study design is the same. The tested methods are the real video and three different SGNify avatars: the solid purple avatar from the first study, the same avatar wearing a black long-sleeved t-shirt, and a fully textured human character adapted from Meshcape [61] (see Fig. 3.9). This study comprises 24 signs (including the four for participant training) from the same ASL dataset [85] but different from those of the first study. Our results reveal that adding the t-shirt and using a fully textured avatar do not benefit actual

Method	Upper Body	Left Hand	Right Hand
FrankMocap [78]	74.93	23.70	19.57
PIXIE [21]	59.09	24.79	20.19
PyMAF-X [94]	68.30	22.51	18.49
SMPLify-SL	55.71	21.14	18.60
SGNify	<b>54.72</b>	<b>20.28</b>	<b>17.44</b>

Table 3.6: Mean TR-V2V error (mm) on fluid sentences.

or perceived sign recognition or the perceived naturalness of the reconstruction. As in the previous study, recognition of SGNify’s reconstructions were not statistically different from real video.

### 3.3.4 SGNify Extensions

We extend SGNify to work on multi-view videos. We used 12 synchronized RGB cameras (see Fig. 3.10) at 90 fps to capture the same participant used in the quantitative evaluation plus two additional signers, a native signer and an interpreter with 17 years of experience. Each participant articulated all signs in our DGS corpus (see Sec. 3.2). A close-up frontal camera is zoomed in to focus on the hands and face of the signer and has a view similar to existing sign-language videos. Another frontal camera captures the whole front body of the participant. Two top-lateral cameras acquire images with a top-down view. Four lateral cameras are placed at hip level and capture the whole body; two are slightly behind the signer, and the other two are slightly in front. Two frontal-lateral cameras also have a full-body view, looking slightly down. Finally, two other frontal cameras, one with a bottom-up view and one with a top-down view, are focused on the hands. The participant stands on a 1.5 m  $\times$  1.5 m platform of adjustable height located in front of a green screen. Then, multi-view SGNify is used to fit SMPL-X. We follow Huang *et al.* [40] to combine the keypoint predictions of each camera using a multi-view consensus approach. A person-specific  $\beta$  is obtained with a 3D scanner. Sample multi-view results are shown in Fig. 3.10 and in the supplemental video at <https://sgnify.is.tue.mpg.de>

Furthermore, we propose a baseline method for continuous SLC (CSLC). CSLC introduces additional challenges, such as the segmentation of sentences into signs; this is an active field of research. When a sentence is given as input, we use Renz *et al.* [76] to segment the input video and then process each segment with SGNify. The first frame of each segment is initialized from the last frame of the previous one.

Besides isolated signs, our corpus contains ten sentences articulated by the three interpreters during the four sessions; one session with a 54-camera Vicon mocap system at 120 fps synchronized with a frontal 4112  $\times$  3008 RGB camera at 60 fps, framing an upper-body view as typically found in SL video (see Sec. 3.2) and three sessions with



Figure 3.10: Multi-view setup with the 12 synchronized RGB cameras and their SGNify reconstruction.



Figure 3.11: Sample frames and reconstructions from segments of the German sentence: *Der Vater muss für die Reparatur seines Autos viel Geld ausgeben.*

the multi-view setup. Each interpreter translated each sentence into DGS, producing different signed versions of the same German sentence.

We conduct an exploratory quantitative study with twelve sentences (ten main sentences and two variations) collected as in Sec. 3.2 and analyzed as in Sec. 3.3.1. Table 3.6 shows the mean TR-V2V error across the twelve sentences for four methods and three body regions. This experiment compares SGNify with FrankMocap [78], PIXIE [21], PyMAF-X [94], and our baseline SMPLify-SL. SGNify achieves the lowest error for the upper body and both hands, beating standard HPE methods. It is interesting to notice that while FrankMocap has a hand-pose error lower than PyMAF-X in our previous quantitative experiment (see Sec. 3.3.1), this is not true in this second experiment. This inconsistency further emphasizes the limitations of a per-frame metric for sign language. In the future, a perceptual study should be conducted to evaluate the recognition of the reconstructed sentences with proficient signers. Such an experiment will give more insights about the next crucial steps for CSLC. Figure 3.11 shows sample frames and SGNify’s reconstructions from a sentence of this exploratory study.

### 3.4 Discussion

Our results show that SGNify performs quantitatively better than standard HPE methods, in particular thanks to the inclusion of our novel linguistic constraints. However, we believe that a per-frame metric is not ideal for SL. To recognize a sign, the temporal evolution is crucial, and this is not captured by V2V. For example, the few slightly inaccurate frames of the SMPLify-SL reconstruction of DOLL (see video on our project page) confused many signers during the perceptual study. Small changes over time can disrupt perception, while overall V2V error remains small. In the end, what matters is whether the meaning is clear to a human. We think a perceptual study provides key insights that complement metric evaluation.

The perceptual study indicated that SGNify significantly outperforms standard HPE methods and, most importantly, produces the first 3D avatars to achieve a sign-recognition accuracy that is not statistically different from the source videos.

### 3.4.1 Limitations and Future Work

SGNify is an optimization-based approach, and as such, it requires substantial computational time ( $\sim 65$  s per frame). Future work should use SGNify’s reconstructions to train a regressor. However, before this step, SGNify can be further improved. Our perceptual study points out the need for enhancements in the face, including facial expressions, tongue and eye movements, mouth morphemes, and eyebrows. While face capture continues to advance as an active field of research, we identified another critical limitation during our evaluation. When examining the avatars from viewpoints other than the frontal view, we observed that signs involving self-touch, *i.e.*, signs where the hands touch each other, the face, or the chest, were not correctly reconstructed, as shown in Fig. 3.12. This error highlights a fundamental challenge in monocular HPE: accurately capturing contact along the camera’s viewing direction. This limitation is not surprising given the inherent difficulty in labeling self-contact from visual data alone. Existing methods are biased toward avoiding contact generation, and capturing contact is particularly challenging in the depth direction of the camera. The ambiguity between contact and close proximity poses a significant obstacle for creating accurate training data at scale. This problem is not unique to SL; all monocular HPE methods encounter similar challenges in detecting self-contact events. To address this fundamental limitation, the following chapters explore bioimpedance sensing as a complementary modality that provides contact information to resolve ambiguities in SLC, and more broadly, in HPE.



Figure 3.12: THANKS (ASL) in the original video and reconstructed by SGNify. The side view reveals the reconstruction error: the hand is positioned significantly far from the body and fails to contact the face, an error not evident from the frontal view.

## Chapter 4

# Detecting Self-Touch Using Bioimpedance

**Note:** This chapter is based on Forte, Vardar, Javot, and Kuchenbecker’s article “Wrist-to-Wrist Bioimpedance Can Reliably Detect Discrete Self-Touch,” which was published in *IEEE Transactions of Instrumentation and Measurement*, Vol. 74, 2025 [28]. Some paragraphs in this dissertation’s Abstract, Introduction, Background, and Discussion are also adapted from this publication.

Our examination of SGNify’s reconstructions reveals that, while linguistic priors significantly improve hand pose estimation, signs involving self-touch along the camera’s viewing axis remain problematic when observed from side views. This limitation is particularly significant because self-touch is integral to SLs; nearly half of the signs involve contact between the left and right hands or between one or both hands and the signer’s face or chest. Their reconstruction is particularly challenging when the hands occlude the touched body part in the camera’s viewing direction. Since visual evidence alone might not be enough in frames where even human observers fail to identify actual contact, we investigate a complementary sensing modality: electrical bioimpedance. This chapter presents a systematic offline study to establish the fundamental capabilities of bioimpedance sensing for self-touch detection. We used prolonged contact poses (~10-second holds) to characterize the bioimpedance response across different frequencies and identify optimal sensing parameters. This controlled approach enables us to establish the sensing principles that will inform real-time dynamic applications in Chapter 5.

While the human sense of touch can reliably detect the occurrence of a self-touch event, the body parts involved, and the surface area of the skin contact, replicating this level of reconstruction remains challenging for sensing devices. Most critically for our SLC goals, reliable self-touch detection could provide pGT contact information to overcome the fundamental limitations identified in monocular pose estimation systems. Beyond SL, the ability to reconstruct self-touch with sensors could positively impact several other fields; it can provide insights into human psychological states and behavior [48], monitor hygiene and virus transmission [50], and enrich wearable or on-skin in-

terfaces [35]. In fact, self-touch actions occur frequently in daily life. These actions are often performed with little or no awareness [64], especially when self-touch is the manifestation of a psychological state, such as stress or task concentration [48]. Facial self-touch is a particularly prevalent human behavior, with individuals unconsciously touching their face approximately 50 times per hour [74], usually with their non-dominant hand [95]. Self-touch also serves communicative purposes, such as partially covering the eyes to express dismay or covering both ears to indicate noise sensitivity, especially in individuals with autism [42].

However, developing reliable self-touch detection systems remains challenging across all application domains. As detailed in Sec. 2.3.1, existing approaches for self-touch detection including vision-based systems [14, 22, 36, 53, 62, 67], biomechanical sensors [64], short-range radar proximity detectors [32], acoustic methods [35, 65], electronic skins [93], and conductivity sensing [96], face significant limitations in terms of accuracy, scalability, or practicality for real-world deployment. Manual annotation approaches [50] are also impractical for large-scale deployment. Among the various sensing modalities, bioimpedance sensing emerges as particularly promising due to its contact specificity, temporal characterization, immunity to visual occlusion, wearability, and power efficiency (see Sec. 2.3.1 for technical principles).

To address the lack of reliable, scalable, and practical self-touch detection systems, our research focuses on a fundamental question: “Which *types of* discrete self-touch poses, *if any*, can bioimpedance-based systems reliably detect across diverse individuals?” We created a dataset consisting of 27 genuine self-touch poses and six adversarial mid-air gestures collected from 30 participants to answer this question. For each pose, we measured the participants’ bioimpedance across a wide range of frequencies by connecting two conductive wristbands to an impedance analyzer, replicating the setup used by Touché [80] to classify five poses. We then conducted a detailed analysis of how the 33 poses of our dataset affect bioimpedance compared to the individual’s baseline. This approach makes our study the first to systematically link self-touch poses to changes in bioimpedance. We also examined the sensitivity of bioimpedance changes to factors such as the body parts involved, the surface area of skin contact, individual characteristics, and external conditions.

Our results show that bioimpedance-based sensing systems hold great promise for reliably detecting skin-to-skin self-touch poses. We identified the specific range of bioimpedance frequencies that are most informative for detecting genuine self-touch only, showing that just the magnitude of bioimpedance at high frequencies needs to be observed to accurately infer these events. Furthermore, we found that bioimpedance changes are strongly influenced by key properties of the contact, such as the body parts involved and the skin contact area. These insights lay the foundation for developing more effective and scalable touch-detection systems that use bioimpedance as the sensing modality, providing researchers with predictive guidelines for which types of self-touch poses are likely to be detectable without extensive additional testing.

## 4.1 Materials and Methods

We conducted an exploratory study to understand how self-touch affects the electrical bioimpedance measured between the left and right wrists in different individuals. We designed this study as a systematic offline characterization to establish the fundamental limits and optimal parameters of bioimpedance-based contact detection. The 10-second pose holds allow us to precisely characterize the response to different contact types. This controlled approach is essential for determining the sensing parameters that will enable real-time dynamic detection in practical applications.

### 4.1.1 Experimental Setup

Two conductive wristbands and a high-quality impedance analyzer (MFIA, Zurich Instruments) were used to precisely measure the participants' bioimpedance, as shown in Fig. 4.1. Specifically, we used a setup similar to the one employed for conducting segmental bioelectrical impedance analysis (BIA) of the upper body, in which alternating current is transmitted from hand to hand through the chest [8]. Mathews and Jovanov [59] recently also identified this wrist-to-wrist configuration as promising for enabling new



Figure 4.1: Experimental setup. The participant wears the wristbands on their left and right wrists. The electrodes in the wristbands are connected to the impedance analyzer through BNC cables. The screen shows the name and a mirrored image of the pose to be mimicked, a message that states whether the participant needs to hold the pose, and their own mirrored video stream captured with an external camera placed above the screen. The laptop computer records the data. The inset shows a top view of both wristbands.

wearable bioimpedance applications. For the electrodes, we explored wet and dry solutions (with different geometries and materials). Since we found consistent results for both, provided good contact was maintained with the skin, we opted for dry electrodes due to their greater practicality. In particular, we used commercial anti-static wristbands (ESD Grounding Wrist Strap, 10 mm Stud) with a conductive circumference of 14 cm (when not stretched) to achieve a skin-to-electrode contact area that is both large and stable. The wristbands were connected to the impedance analyzer through shielded cables to reduce noise.

The impedance analyzer was set to operate in four-terminal measurement mode to remove the resistive effect of the wires. We used a parallel resistance and capacitance as the equivalent circuit and a high-accuracy sweep [99]. We used electrical bioimpedance spectroscopy (EBIS) to measure the wrist-to-wrist impedance of the participant from 100 Hz to 5.1 MHz over 100 logarithmically spaced frequencies. This frequency range was chosen to include and extend beyond all frequencies commonly used for BIA. Single-frequency BIA typically measures the phase angle of the bioimpedance at only 50 kHz [51]. Multi-frequency BIA, the most widespread and well-known application of EBIS [7] used for assessment of total body composition and fluid distribution, sweeps from 5 kHz to 200 kHz because this range has the highest reproducibility, even though day-to-day coefficients of variation increase for frequencies below 50 kHz [51]. Nonetheless, our goal is different from estimating body composition, and we thus measured the bioimpedance across the full range offered by our impedance analyzer.

The experiments were conducted in a temperature-controlled research laboratory. During the experiment, the participant sat in front of a screen displaying a MATLAB graphical user interface (GUI) with a mirrored image of the pose to mimic, their mirrored video stream to check whether they were performing the pose correctly, and a message that told them whether to hold or release the pose. The experimenter used a nearby laptop computer to record the impedance measurements and the participant's screen.

### 4.1.2 Experimental Protocol

We recruited a total of 30 participants, aged  $31.2 \pm 6.4$  years (mean  $\pm$  standard deviation). Of these, 15 self-reported to be male and 15 to be female, with six participants being left-handed and the remainder right-handed. The participants represented diverse ethnic backgrounds, including Asian, White, and Hispanic, and had a BMI of  $22.9 \pm 3.0$  kg/m<sup>2</sup>. None of them had current or past sensory-motor disabilities. Approval of the experimental procedure for this study was granted by the Ethics Council of the Max Planck Society under the Haptic Intelligence Department's framework agreement (protocol number F013C). All participants provided informed consent to participate in the study before data collection. People not employed by our organization were offered a nominal hourly payment.

At the start of the study, the experimenter introduced the 33 poses and the baseline no-touch pose depicted in Fig. 4.2. As previously stated, self-touch refers to using one's

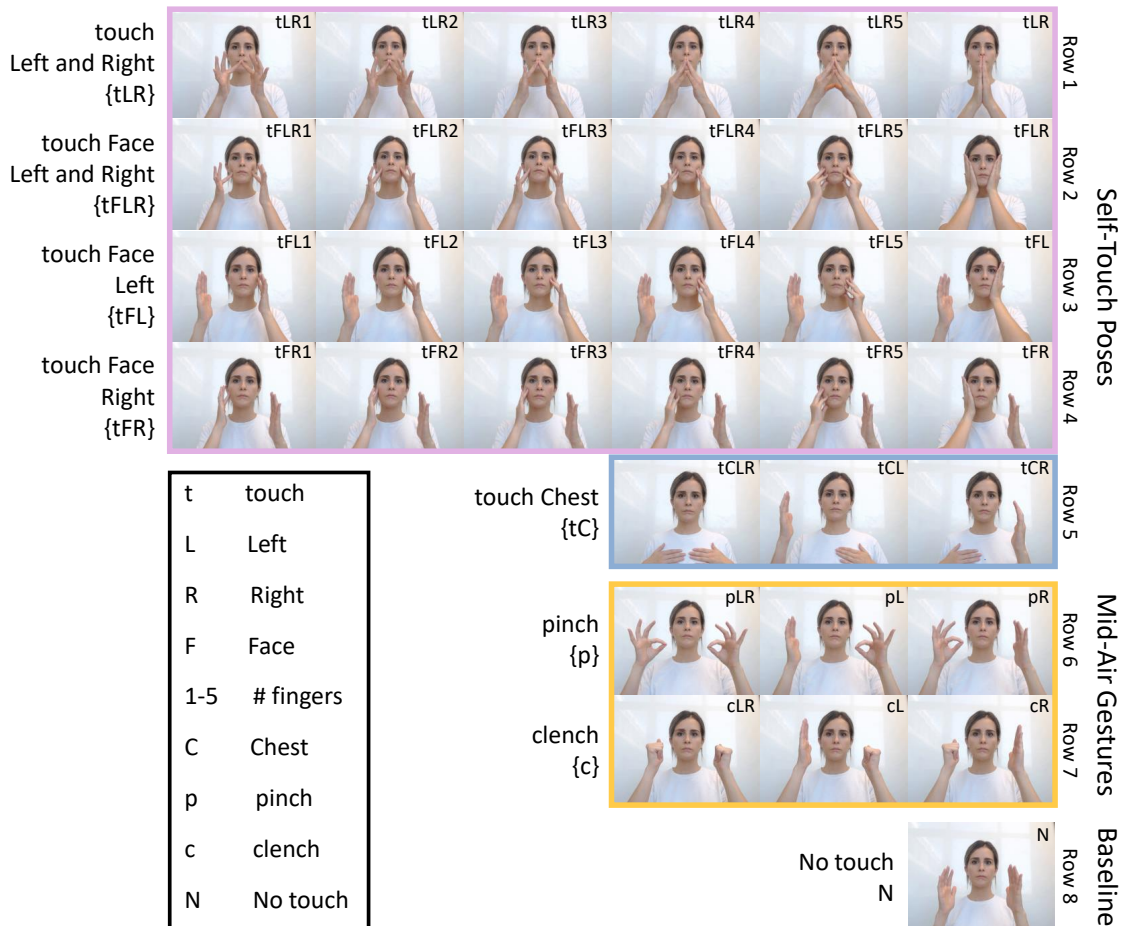


Figure 4.2: Poses performed during the experiment. Each row shows a group of related poses, as labeled at the left. The first five rows show self-touch poses; they involve touch interactions between the left and right hands {tLR}, contact with the face by either both hands {tFLR}, the left hand {tFL}, or the right hand {tFR}, and touching the chest with one or both hands {tC}. Skin-to-skin and skin-to-clothing self-touch poses are highlighted in pink and blue, respectively. Rows six and seven show adversarial mid-air gestures: pinching {p} and clenching {c}, performed with either one or both hands. The last row depicts the no-touch condition N.

hand(s) to contact another body part, either on the skin or through clothing, with discrete facial self-touches being particularly common. These poses are particularly relevant for SL applications, where hand-to-hand, hand-to-face, and hand-to-body contacts are integral to communication [88]. We thus chose 27 poses (rows 1–5) that involve various levels of skin-to-skin contact (rows 1–4) and skin-to-clothing contact (row 5). Finally, we added three variations of two mid-air hand gestures (rows 6–7) that are commonly used in human-computer-interaction applications [39]: pinching the thumb and index fingers together and clenching all of the fingers into a fist. These adversarial gestures can be performed with either one or both hands. Since they involve within-hand skin-to-skin contact, they might be erroneously detected as self-touch poses by bioimpedance-based self-touch systems. The 34<sup>th</sup> pose is the no-touch pose, which gives us the baseline for the user’s wrist-to-wrist bioimpedance (N; row 8). To better understand the underlying physiological mechanisms, we purposefully selected some poses that have electrical topologies similar to one another (*e.g.*, performing the same action with the left or right hand or with a different number of fingers). For conciseness, the actions of touch, pinch, and clench are labeled with lowercase letters (t, p, and c). Capital letters denote the body parts in contact, *i.e.*, left hand, right hand, face, and chest (L, R, F, and C). A numeral at the end of a touch label indicates the number of fingertips in contact (1, 2, 3, 4, or 5). We standardized the fingers used: 1 means only the index finger; 2 means index and middle; 3 means index, middle, and ring; 4 means index, middle, ring, and pinky; and 5 means all five fingers. When there is no number, the entire hand makes contact, including the full fingers and the palm.

After the experimenter’s introduction, the participant practiced the poses to ensure comprehension. They then completed a short training to become familiar with the experimental procedure before data collection. The experiment required the participant to complete three measurement cycles; within each cycle, the GUI presented the 33 poses in random order, and each pose was preceded by the no-touch baseline pose. Recording the baseline bioimpedance before each pose, instead of only once within each cycle, is fundamental to our approach, since the user’s bioimpedance changes continuously over time due to internal [31] and external factors (such as changes in room humidity, temperature, and wristband position). Each sweep took about 10 seconds, and the participant held the pose slightly before and after the sweep. After each pair of bioimpedance measurements, *i.e.*, baseline and pose, the participant could take a short break, and they had a longer break between cycles.

### 4.1.3 Dataset

The raw dataset includes 1,188,000 data points, which correspond to the magnitude and phase measurements at the 100 frequencies for the 198 trials of each of the 30 participants. In addition, the dataset contains the participants’ individual characteristics (*i.e.*, sex, handedness, ethnicity, and BMI), the material of the clothing touched during the chest contacts, and their chosen starting position for the poses (*i.e.*, elbows on the table

or in the air). We did not include the age of the participants in the dataset due to its low variability.

The variable indicating the performed pose, *i.e.*, PoseID, has 33 levels (24 skin-to-skin contacts, three skin-to-clothing contacts, and six adversarial mid-air gestures). We refer to all other variables as attributes. Among the attributes, sex, handedness, and starting position have only two levels each. In the case of double ethnicity, we considered the primary one; this choice led to three ethnicity levels: Asian (14 participants), White (14 participants), and Hispanic (two participants). We used the WHO categorization [91] to label the BMI as underweight, normal weight, pre-obesity, and obesity (grouping the three obesity classes together); our participants' BMIs span this full spectrum: one underweight, 21 normal weight, seven pre-obesity, and one obese. Finally, we categorized the clothing materials into three groups. The clothing of six participants did not have any label, so we classified their materials as Unknown. The other two groups were formed based on the clothing conductivity [37]: low-conductivity synthetic polymers, cotton, and linen (24 participants combined) were separated from clothing containing wool (two participants), which is generally more conductive.

After collecting the dataset, we created the new variable  $\Delta\text{Bioimpedance}$  that represents the difference between each baseline measurement (the no-touch pose N) and the pose that immediately followed it (*i.e.*,  $\Delta\text{Bioimpedance} = \text{Baseline} - \text{Pose}$ ). Over all experimental sessions, ten bioimpedance measurements lacked either the magnitude or the phase at the highest frequency of 5.1 MHz, presumably due to a data-saving error. These ten incomplete data points were thus removed from the dataset, leaving 593,990 magnitude values and 593,990 phase values.

#### 4.1.4 Statistical Analysis

We conducted two sets of statistical analyses to understand which types of discrete self-touch poses, if any, bioimpedance-based systems can reliably detect across diverse individuals.

##### Detection of Self-Touch Poses

First, to understand the link between genuine self-touch poses and variations in bioimpedance across individuals, we analyzed the collected data using the following mixed-effects model, specified in R formula notation<sup>5</sup>:

$$\Delta\text{Bioimpedance} \sim -1 + \text{PoseID} + \text{Sex} + \text{Handedness} + \text{Ethnicity} + \text{BMI} + \text{ClothingMaterial} + \text{StartPosition} + (1 \mid \text{ParticipantID}),$$

<sup>5</sup>In R formula notation:  $\sim$  means “is modeled as,”  $-1$  provides separate estimates for each level instead of comparing everything to one baseline,  $+$  means “include this variable,” and  $(1 \mid \text{variable})$  accounts for the fact that measurements from the same unit are correlated.

where  $\Delta\text{Bioimpedance}$  is the dependent variable, and  $\text{PoseID}$  is the independent variable. The impact of individual characteristics (*i.e.*, sex, handedness, ethnicity, and BMI) and external factors (*i.e.*, clothing material and starting position) is analyzed by including them as fixed effects. The ID of the participant is modeled as a random effect to account for within-subject variability, as each participant performed the same poses three times. The model was separately fit for the magnitude and the phase using the ‘lme4’ package in R. Concerns about any potential violations of model assumptions should be alleviated by the robustness of mixed-effects models [81]. After confirming the overall significance of each fixed effect with a Type II Wald chi-square test, we performed a False Discovery Rate post-hoc correction; this particular correction was chosen due to the exploratory nature of this analysis [30].

We deepened our understanding of the influencing factors through a follow-up sensitivity analysis using a leave-one-variable-out approach. In this analysis, we systematically removed one fixed effect at a time, refit the model, and evaluated how the exclusion of each variable impacted the model’s performance by examining changes in the conditional and marginal  $R^2$  value, root mean square error (RMSE), Akaike information criterion (AIC), and Bayesian information criterion (BIC).

### **Influence of Body Parts and Skin Contact Area on Skin-to-Skin Bioimpedance Changes**

Due to their prevalence and relevance for several application fields, our second analysis focused on the skin-to-skin self-touch poses (rows 1–4 in Fig. 4.2), aiming to understand how the body parts and skin contact area involved in these contacts influence bioimpedance changes. The 24 poses tested can be divided into either four groups based on the body parts (each group is a row from rows 1–4 in Fig. 4.2) or six groups based on the contact size (each group is a column from rows 1–4 in Fig. 4.2). We thus introduced the respective discrete variables  $\text{TopologyID}$  and  $\text{ContactAreaID}$  and separately fit the following mixed-effects model to the magnitude and phase of the bioimpedance:

$$\Delta\text{Bioimpedance} \sim -1 + \text{TopologyID} * \text{ContactAreaID} + \text{TopologyID} + \text{ContactAreaID} + \text{Sex} + \text{Handedness} + \text{Ethnicity} + \text{BMI} + \text{ClothingMaterial} + \text{StartPosition} + (1 \mid \text{ParticipantID}).$$

We then used a Type II Wald chi-square test to check the significance of the interaction effect of  $\text{TopologyID}$  and  $\text{ContactAreaID}$ . When the interaction effect was not significant, we analyzed  $\text{TopologyID}$  and  $\text{ContactAreaID}$  independently. When significant, we explored the simple effects of each of the two variables at the different levels of the other variable. After confirming the overall significance of each fixed effect with a Type II Wald chi-square test, we corrected the significance using Tukey’s HSD for the pairwise comparisons of  $\text{TopologyID}$  and  $\text{ContactAreaID}$  and Bonferroni for the attributes.

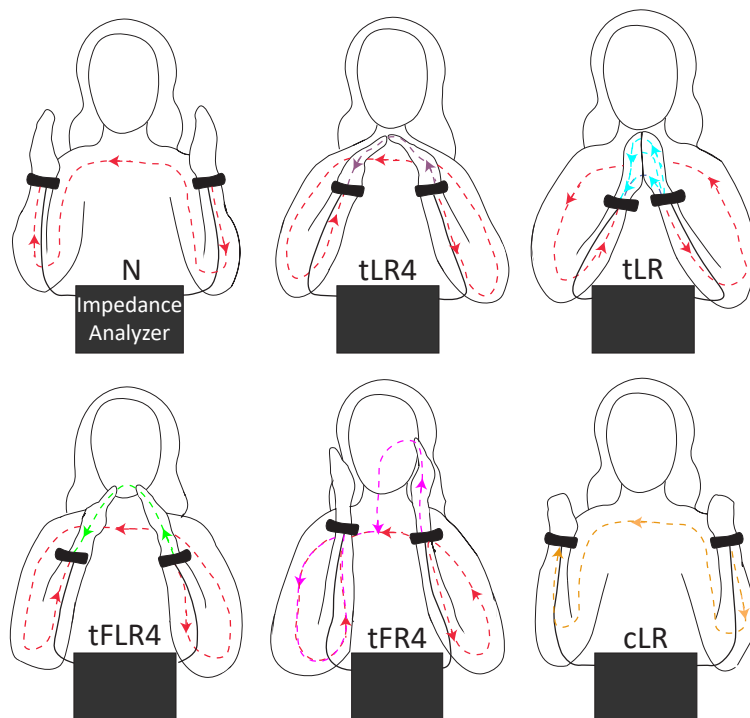


Figure 4.3: Mirrored illustrations of current pathways for six sample poses. The red pathway through the user's arms and across their shoulders is always present. Touching the hands together or to the face gives the current a second way to flow from wrist to wrist. Performing a hand gesture like clenching does not alter the circuit topology but might have a somewhat different impedance, shown in yellow.

## 4.2 Results

Figure 4.3 shows schematics for the baseline no-touch pose (N) and five other selected poses to illustrate their different conductivity topologies. This set includes four skin-to-skin self-touch poses, *i.e.*, four fingertips of both the left and right hands touching each other (tLR4), the entire hands touching each other (tLR), four fingertips of both the left and right hands touching the face (tFLR4), and four fingertips of only the right hand touching the face (tFR4), as well as one mid-air gesture, *i.e.*, clenching both hands (cLR). For all six of these poses, the three wrist-to-wrist bioimpedance measurements (magnitude and phase across frequencies) from a sample participant are depicted in Fig. 4.4. The other poses follow similar trends but are omitted for presentation clarity. Both the magnitude and the phase of the measured bioimpedance show systematic differences as a function of frequency. The magnitude decreases with increasing frequency until a change in trend occurs at the highest frequencies, where it rises due to the dominant influence of the parasitic inductance from the cables. In general, it is possible to notice substantial differences between the bioimpedance magnitudes across poses at high frequencies (200 kHz to 4.5 MHz, see insets in Fig. 4.4). The baselines are also relatively stable in this range, hinting that  $\Delta$ Bioimpedance will depend mainly on the pose itself. The bioimpedance phase, instead, seems to differ across poses mainly at the low and middle frequencies, though the baselines also fluctuate greatly in this range.

### 4.2.1 Detection of Self-Touch Poses

Figure 4.5 reports the number of fixed effects that caused a significant change in magnitude (Fig. 4.5a) or phase (Fig. 4.5b) from the preceding baseline. The PoseID variable is represented by the bars in pink, blue, and yellow (skin-to-skin contacts, skin-to-clothing contacts, and adversarial mid-air gestures, respectively). The gray bars represent instead the attributes of sex, handedness, ethnicity, BMI, clothing material, and starting position. Lighter color shades indicate significance levels of  $p < 0.05$ , while darker shades mark  $p < 0.001$ . We start by analyzing the variable PoseID and then focus on the attributes. This subsection concludes with the results of the associated sensitivity analysis.

It is important to note that when referring to our results, we use the term “detection” in a statistical sense rather than in the context of machine-learning classification: we refer to significant deviations in bioimpedance from the baseline level.

#### PoseID

As seen in Fig. 4.5a, starting at frequency 58 (51.3 kHz), all skin-to-skin poses exhibited a bioimpedance magnitude that significantly differed from the preceding baseline ( $p < 0.05$ ). In particular, from approximately 137.5 kHz to 4.1 MHz (from frequency 67 to frequency 98), they were all detected with  $p < 0.001$ . Starting from frequency 80 (570.9 kHz), all skin-to-clothing contacts also caused significant changes in

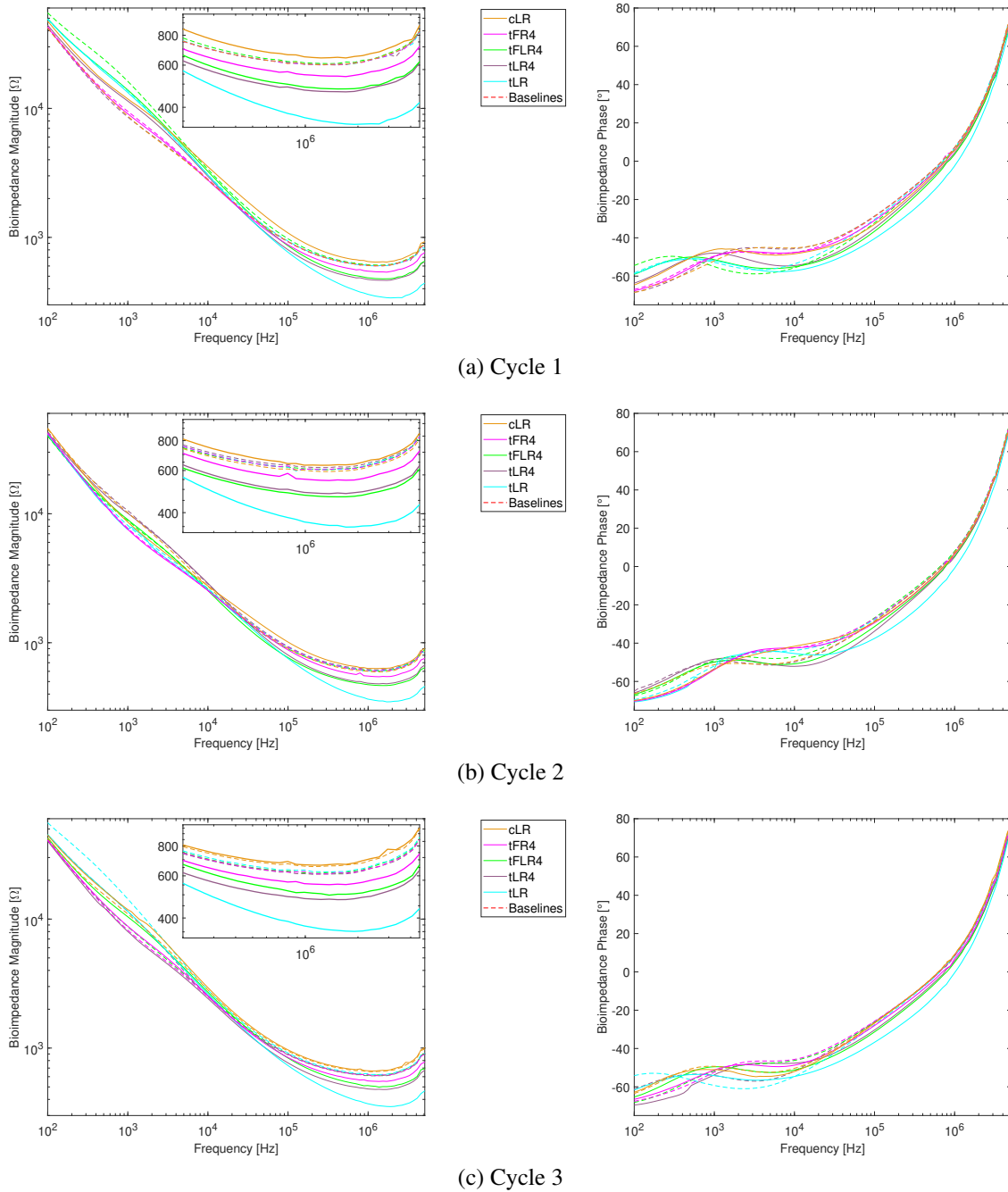


Figure 4.4: Wrist-to-wrist bioimpedance of the three cycles of participant 1 for five selected poses (solid lines), i.e., tLR4, tLR, tFLR4, tFR4, cLR, and their respective baselines (dashed lines). The magnitude (left) and phase (right) were measured for (a) Cycle 1, (b) Cycle 2, and (c) Cycle 3 from 100 Hz to 5.1 MHz. The insets in the magnitude plots show the zoomed-in responses from 200 kHz to 4.5 MHz, which are relatively stable across cycles.

the bioimpedance ( $p < 0.05$ ), and the significance was  $p < 0.001$  at many frequencies for touches performed with two hands (tCLR). When only one hand was used, we observed more detections for the left hand (tCL) than for the right hand (tCR). Among the adversarial mid-air gestures, clenching significantly differed from N ( $p < 0.05$ ) from 100 Hz to 213.1 kHz. Notably, when performed with both hands (cLR) or with the left hand alone (cL), clenching demonstrated substantial significance ( $p < 0.001$ ) across multiple frequencies, with a higher rate of significance observed in the two-hand condition. No pinching pose was detected at any frequency.

In Fig. 4.5b, it can be seen that most poses showed a significant deviation in the bioimpedance phase around 50 kHz. In this region, the only skin-to-skin contact that exhibited no significant difference from the baseline was the use of a single fingertip to touch the face (tFR1). All skin-to-clothing poses were detected ( $p < 0.001$ ) starting from 1.4 MHz. Similar to before, among the adversarial mid-air gestures, only a few clench poses were detected (cLR, followed by cR and cL) at low and medium frequencies. Starting from frequency 95 (2.9 MHz), the results became less consistent due to the increased impact of the parasitic inductance on the phase.

### Attributes

Sex emerged as the only significant attribute ( $p < 0.05$ ) at some low and moderately high frequencies in the magnitude, and had a significant impact ( $p < 0.05$ ) also at low frequen-

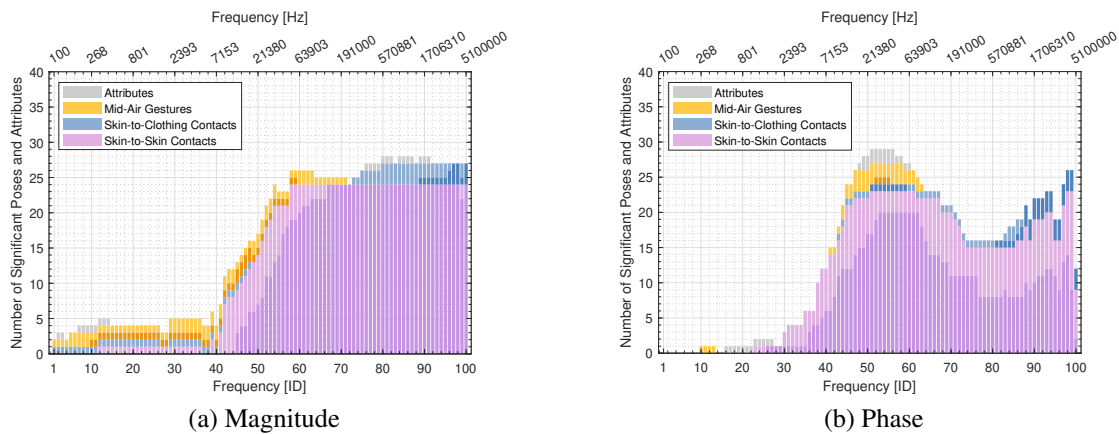


Figure 4.5: Stacked bar chart illustrating the number of significant effects on either (a) magnitude or (b) phase of the measured bioimpedance change. Pink refers to the skin-to-skin contacts (maximum 24), blue to the skin-to-clothing contacts (maximum three), and yellow to the mid-air gestures (maximum six). Gray combines all the other fixed effects, i.e., the attributes of sex, handedness, ethnicity, BMI, clothing material, and starting position. The lighter colors show significant differences at  $p < 0.05$ , and the darker colors mark  $p < 0.001$ .

cies in the phase. Ethnicity and clothing material significantly affected the bioimpedance phase ( $p < 0.05$ ) in the medium frequencies. Handedness, BMI, and the starting position did not influence the measurements.

## Sensitivity Analysis

Based on these results, we conducted a sensitivity analysis on only the magnitude and averaged the results within three frequency clusters: low (from 100 Hz to 7.2 kHz), medium (from 8.0 kHz to 213.1 kHz), and high (from 237.8 kHz to 5.1 MHz). Overall, the model with all fixed effects consistently performed best across all metrics, in particular at high frequencies. As expected, PoseID was the most critical predictor; removing it caused dramatic decreases in both conditional and marginal  $R^2$  values (89–97%) and substantial increases in RMSE, AIC, and BIC. In contrast, removing the other variables led to minor (Ethnicity, BMI, and ClothingMaterial) or minimal (Sex, Handedness, StartPosition) changes in the metrics.

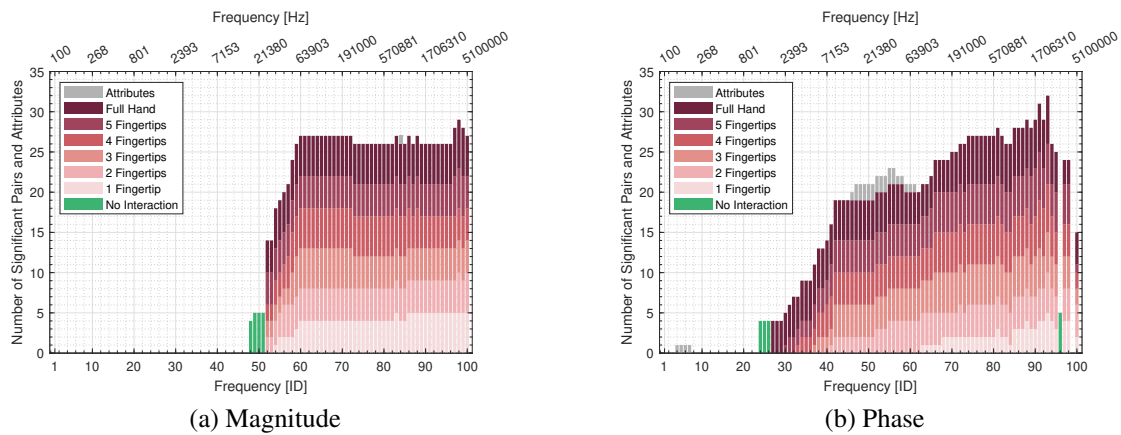


Figure 4.6: Stacked bar chart illustrating the number of significant ( $p < 0.05$ ) TopologyID pairwise comparisons and effects for either (a) magnitude or (b) phase of the bioimpedance change (maximum six for each label). The green bars (no interaction) mark the pairwise comparisons that are independent of the ContactAreaID, i.e., when the interaction effect between ContactAreaID and TopologyID is not significant. The other six labels refer to each of the six levels of ContactAreaID, i.e., when the interaction effect between ContactAreaID and TopologyID is significant. Gray combines all the other fixed effects, i.e., the attributes of sex, handedness, ethnicity, BMI, clothing material, and starting position.

## 4.2.2 Influence of Body Parts and Skin Contact Area on Skin-to-Skin Bioimpedance Changes

We streamline the presentation of our second analysis by reporting results only for  $p < 0.05$ . We first report the impact of TopologyID and ContactAreaID on the dependent variable  $\Delta$ Bioimpedance and then conclude the subsection focusing on the attributes (gray bars).

### TopologyID and ContactAreaID

There was no significant interaction effect between TopologyID and ContactAreaID in the magnitude from frequency 1 (100 Hz) to frequency 51 (23.9 kHz), and in the phase from frequency 1 (100 Hz) until frequency 26 (1.5 kHz), as well as at frequencies 96 (3.3 MHz) and 99 (4.6 MHz).

Figure 4.6 shows the results of TopologyID. Until frequency 48 (17.2 kHz), none of the six pairwise comparisons between the four touch topologies ( $\{tFR\}$ ,  $\{tFL\}$ ,  $\{tFLR\}$ , and  $\{tLR\}$ ) showed significant differences in how they changed the bioimpedance magnitude (Fig. 4.6a). From this frequency to frequency 51 (23.9 kHz), most pairwise comparisons were significant, except for the left- versus right-hand face touch ( $\{tFL\}$  versus  $\{tFR\}$ ), which was not significant throughout this range, and between the hands touching directly or through the face ( $\{tFL\}$  versus  $\{tFLR\}$ ), which was non-significant only at frequency 48. Starting at frequency 52, overall, the pairs comparing the use of one hand with two hands ( $\{tFL\}$  or  $\{tFR\}$  versus  $\{tLR\}$  or  $\{tFLR\}$ ) showed significant differences for most frequencies and contact sizes; the less-significant pairwise comparisons were again  $\{tFL\}$  versus  $\{tFR\}$  followed by  $\{tLR\}$  versus  $\{tFLR\}$ . Furthermore, the contact area with the highest number of significant pairwise comparisons was the full hand.

Regarding the phase, as shown in Fig. 4.6b, until frequency 24 (1.2 kHz) and at frequency 99 (4.6 MHz), none of the six pairwise comparisons was significant. From frequency 24 (1.2 kHz) to frequency 26 (1.5 kHz), the only non-significant comparisons were again  $\{tFL\}$  versus  $\{tFR\}$  and  $\{tLR\}$  versus  $\{tFLR\}$ , with the latter being the only non-significant comparison at frequency 96 (3.3 MHz). Starting from frequency 27 (1.7 kHz), overall, the number of significant pairwise comparisons increased proportionally with the frequency and the contact size. At a few high frequencies, all pairwise comparisons were significant when either four or five fingertips were used. With the full hand(s), all comparisons were significant for most of the high frequencies. Starting with frequency 95 (2.9 MHz), the results became less consistent, similar to the behavior witnessed in the first analysis.

Figure 4.7 shows the results of ContactAreaID. As the frequency increases, so does the difference in the changes in bioimpedance magnitude between contact sizes. At lower frequencies, only a few pairs with large contact-area differences, *e.g.*, one fingertip versus the full hand, were significantly different. However, at high frequencies all 15 pairwise comparisons were significant when both hands were used, *i.e.*,  $\{tLR\}$  and  $\{tFLR\}$ , and

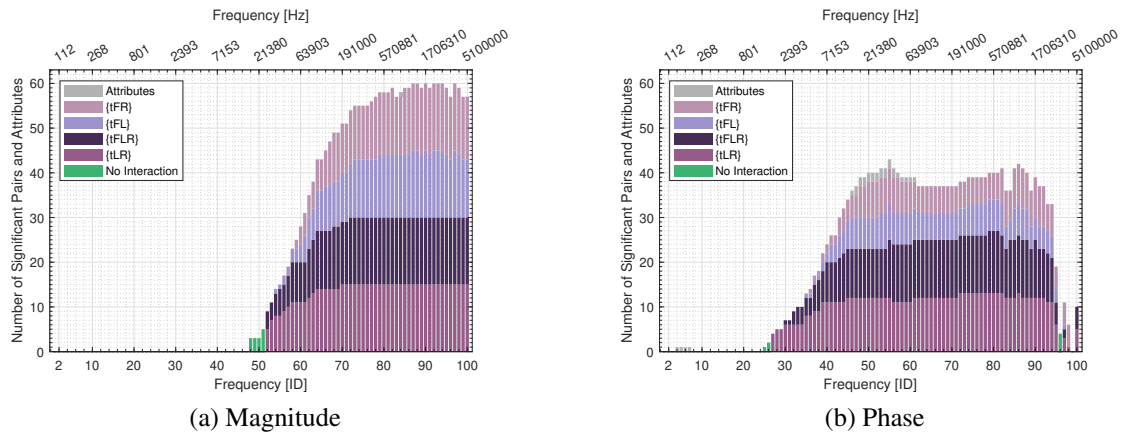


Figure 4.7: Stacked bar chart illustrating the number of significant ( $p < 0.05$ ) ContactAreaID pairwise comparisons and effects for either (a) magnitude or (b) phase of the bioimpedance change (maximum 15 for each label). The green bar (no interaction) defines the pairwise comparisons independently of the TopologyID, i.e., when the interaction effect between TopologyID and ContactAreaID is not significant. The other four labels refer to each of the four levels of TopologyID, i.e., when the interaction effect between TopologyID and ContactAreaID is significant. Gray combines all the other fixed effects, i.e., the attributes of sex, handedness, ethnicity, BMI, clothing material, and starting position.

the large majority when only one hand was used, i.e., {tFL} and {tFR}, with some frequencies where all comparisons were significant. The phase followed a similar trend but resulted in fewer significant comparisons, with the number of significant comparisons reducing by about 50% when only one hand was used.

### Attributes

In the magnitude, only one attribute emerged as significant, and it affected only one frequency: sex at 884.6 kHz. For the phase, the starting position had a significant impact at some very low frequencies, while ethnicity, either alone or with clothing, had a significant impact at the medium frequencies.

## 4.3 Discussion

This study explored the feasibility of using wrist-to-wrist bioimpedance to detect discrete self-touch poses, focusing on both skin-to-skin and skin-to-clothing interactions. Our findings indicated that skin-to-skin self-touch poses caused significant changes in bioimpedance magnitude at high frequencies (between 237.8 kHz and 4.1 MHz), making

them consistently detectable across individuals. Skin-to-clothing contacts, while discernible, presented more challenges due to the lower conductivity and high variability of clothing. Importantly, adversarial mid-air gestures did not cause significant deviations in the bioimpedance magnitude at high frequencies, which is crucial for distinguishing genuine self-touch from within-hand contacts.

We attribute our measurement method's higher success at detecting self-touch poses, in particular skin-to-skin contact, and its lower responsiveness to mid-air hand gestures to these poses' *different circuit topologies*. As visible in Fig. 2.2 (in Sec. 2.3.1) and Fig. 4.3, contact between the hands (or with the face or clothing) creates an additional current pathway parallel to the baseline circuit, which decreases the point-to-point bioimpedance magnitude (compare dashed and solid lines at high frequencies in Fig. 4.4). Conversely, mid-air hand gestures, such as clenching or pinching the fingers, do not create a new current pathway between the wrists, therefore showing bioimpedance magnitude values more similar to baseline. Between these mid-air gestures, clenching was likely more detectable due to its increased muscle contraction and skin stretching, two mechanisms that influence bioimpedance measurements [3, 12]. The phase of the bioimpedance detected the highest number of poses with the lowest p-values around 50 kHz (in line with prior work on BIA). However, in this range of frequencies, all clenching poses were also detected, and the clothing material had a significant impact. These results make phase less appealing than magnitude for use in detecting self-touch.

Moreover, we showed that the body parts and the contact area involved in skin-to-skin contacts played a critical role in the bioimpedance magnitude response. The use of one or two hands showed significant differences: in one-handed self-touches, the current has to travel through a longer parallel electrical pathway, which goes from the fingers down through the neck and out the arm, whereas in two-handed self-touches, the shortest parallel pathway goes from one hand to the other either directly or through the face (Fig. 4.3). Furthermore, touching the hands directly or through the face differed significantly at some frequencies, despite the similar pathway lengths and contact areas of these two topologies (see Fig. 4.3). The variability across frequencies is likely due to the softer nature of the cheeks, which led to higher variations in the contact area and applied force compared to the hand-to-hand contacts. Finally, when considering the one-handed self-touches, touching the face with the right or left hand creates a mirrored circuit topology, which could explain their very similar bioimpedance changes. Interestingly, at high frequencies and with large contact areas, the phase showed differences between the usage of the left or right hand, with the right hand showing higher phase values independent of the handedness of the participant.

For the contact size, we observed a systematic decrease in the bioimpedance magnitude with increasing contact size (compare tLR4 and tLR in Fig. 4.4), as larger contact areas facilitate electrical flow. Furthermore, at high frequencies, the only non-significant pairs were in the one-handed facial self-touch poses. Besides the variations due to the softness of the cheeks, the longer path of these poses, compared to the two-handed ones, results in a higher total impedance and, consequently, a smaller deviation from the base-

line impedance (see the pink and green lines in Fig. 4.4), which in turn decreases the relative differences between contact areas. In the phase, the difference between the presence or absence of the face and the use of one or two hands was emphasized even more. Given these results, we believe that subtle differences, such as the use of one or two fingers during one-handed facial contact, may become distinguishable when analyzing the full bioimpedance spectrum rather than examining each frequency independently, closer to the approach used in Touché for classifying poses [80].

Regarding the attributes, our analyses showed that handedness and BMI did not have any significant impact. Instead, sex, ethnicity, clothing, and starting position showed frequency-dependent effects: sex influenced magnitude at low and high frequencies and phase at low frequencies; men are generally larger than women, which will tend to give them higher wrist-to-wrist bioimpedance. Ethnicity and clothing had an impact on the phase at medium frequencies. Phase was also sensitive to the starting position. We believe that the impact of sex at high frequencies, which emerged as the best range to detect only self-touch poses, does not undermine the generalizability of the tested approach. First, our sensitivity analysis indicated that the impact of the sex attribute is minor. Second, this information could easily be provided as an input to the sensing system and is constant for a user, unlike clothing and starting position, whose influence on the phase makes phase even less attractive than magnitude for detecting self-touch.

Our findings demonstrate that bioimpedance-based sensing systems can effectively detect skin-to-skin self-touch poses across a diverse range of individuals. This approach performs particularly well when both hands are in contact (whether directly or through another body part) and is especially effective with larger contact areas. Based on our understanding of the underlying physiological mechanisms, we believe this method's ability to detect self-touch is independent of factors like fine-grained contact location, hand pose, and the specific fingers involved, meaning bioimpedance magnitude at high frequencies can be used to detect skin-to-skin self-touch beyond the chosen set of poses. Our controlled study design with prolonged contact poses was essential for characterizing bioimpedance responses and determining optimal sensing parameters, *i.e.*, magnitude only at frequencies between 237.8 kHz and 4.1 MHz. While this study did not capture the temporal dynamics of natural movements, it provided the foundational understanding needed to design real-time systems. We publicly share the dataset and code associated with this paper to allow further investigations of these phenomena [26].

### 4.3.1 Limitations and Future Work

We observed that one of the primary limitations of this sensing method is its lower sensitivity (*i.e.*, lower significance level) in detecting skin-to-clothing contacts compared to skin-to-skin contacts. Unlike skin-to-skin contacts, which create direct conductive pathways, skin-to-clothing interactions introduce complex electrical interfaces, increasing variability in bioimpedance readings. Additionally, the diversity of fabrics affects conductivity, leading to inconsistent deviations from the baseline pose. Based on our find-

ings, we believe that other sensing modalities may be more effective for skin-to-clothing detection. However, since our dataset included only three skin-to-clothing poses, reflecting their lower frequency in SL scenarios, a more thorough analysis is necessary before discarding the use of bioimpedance sensing for these contacts.

Furthermore, our experiments were conducted under controlled laboratory conditions, allowing us to systematically explore the capabilities of bioimpedance sensing. However, this controlled environment does not reflect the complexities of dynamic real-world conditions, where external factors, such as drastic temperature fluctuations, may influence both the fabric's conductivity and the changes in the user's bioimpedance. To bridge this gap between laboratory validation and practical deployment, future research should focus on testing the studied measurement approach in more naturalistic settings. Additionally, future studies should include participants from a broader range of age groups to better understand the potential influence of age on bioimpedance changes.

Despite these methodological considerations, this controlled study establishes the potential of bioimpedance as a sensing method for self-touch detection. Given the prevalence of hand-to-hand and hand-to-face contacts in SL, this technology can play a crucial role in detecting self-touch in SLC, as well as in other HPE areas. Such an application is explored in the next chapter, which introduces a framework that integrates bioimpedance-based contact detection with vision-based HPE to improve pose estimates across various applications, including SL. Beyond HPE applications, the broader implications of reliable self-touch detection extend to multiple domains. Such a system could also enhance healthcare by monitoring face-touching behaviors in hospitals to mitigate infection risks, especially during outbreaks like COVID-19. In behavioral and psychological research, it could track stress-induced self-touch patterns, providing valuable data for emotional assessments and therapeutic interventions. Additionally, it could improve human-computer interaction by distinguishing between contact and hovering in virtual-reality environments. In applications where users interact with the external environment, however, contact with other people or conductive objects and surfaces would generate bioimpedance changes similar to those of self-touch events. Disambiguating these cases will require integrating complementary sensing modalities or contextual information. Finally, to realize these diverse applications and move beyond laboratory settings, practical implementation requires a portable and wearable system. This technological challenge is addressed in the next chapter, which also presents a miniaturized version of our research-grade impedance analyzer.

# Chapter 5

## Integrating Bioimpedance Sensing into Human Pose Estimation

**Note:** This chapter is based on Forte, Athanasious, Ballardini, Bartels, Kuchenbecker, and Black’s article “Contact-Aware Refinement of Human Pose Pseudo-Ground Truth via Bioimpedance Sensing,” which was published in *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV), 2025, pp. 5071-5080* [25] as well as on a related provisional patent filing concerning the miniaturized sensor [27]. Some paragraphs in this dissertation’s Abstract, Introduction, Background, and Discussion are also adapted from these works.

While the examination of SGNify’s results from a side view (Fig. 3.12) highlights the challenge of reconstructing self-touch, the previous chapter shows that bioimpedance sensing can be used to directly measure the presence of contact, in particular in skin-to-skin contacts. This chapter thus presents a practical system that integrates this complementary modality in HPE methods: BioTOUCH (Bioimpedance Timing for Understanding Contact in Humans). BioTOUCH uses bioimpedance sensing for temporal detection, *i.e.*, when contact occurs, and computer vision for spatial localization, *i.e.*, where the contact occurs, creating a complete solution for contact-aware HPE. This work focuses on self-touch due to its prevalence in SL. However, BioTOUCH can be extended to other self-contact scenarios with appropriate electrode placement, addressing the difficulty of generating reliable pGT data for self-contact interactions.

The failure of standard HPE methods stems largely from the inherent difficulty of obtaining reliable training data with accurate contact labels. Existing approaches have attempted various solutions to address this challenge. SMPLify-XMC [67] generates pGT for the Mimic-The-Pose (MTP) dataset, a dataset where volunteers mimic reference poses of the 3D Contact Poses (3DCP) dataset [67] while being photographed; however, the quality of the resulting dataset depends on how accurately people imitate the depicted poses. Furthermore, 3DCP’s contact detection relies on geometric thresholds combining Euclidean and geodesic distances, thus risking confusing close proximity with contact. Similarly, the pGT from the Discrete Self-Contact (DSC) dataset [67] depends on the

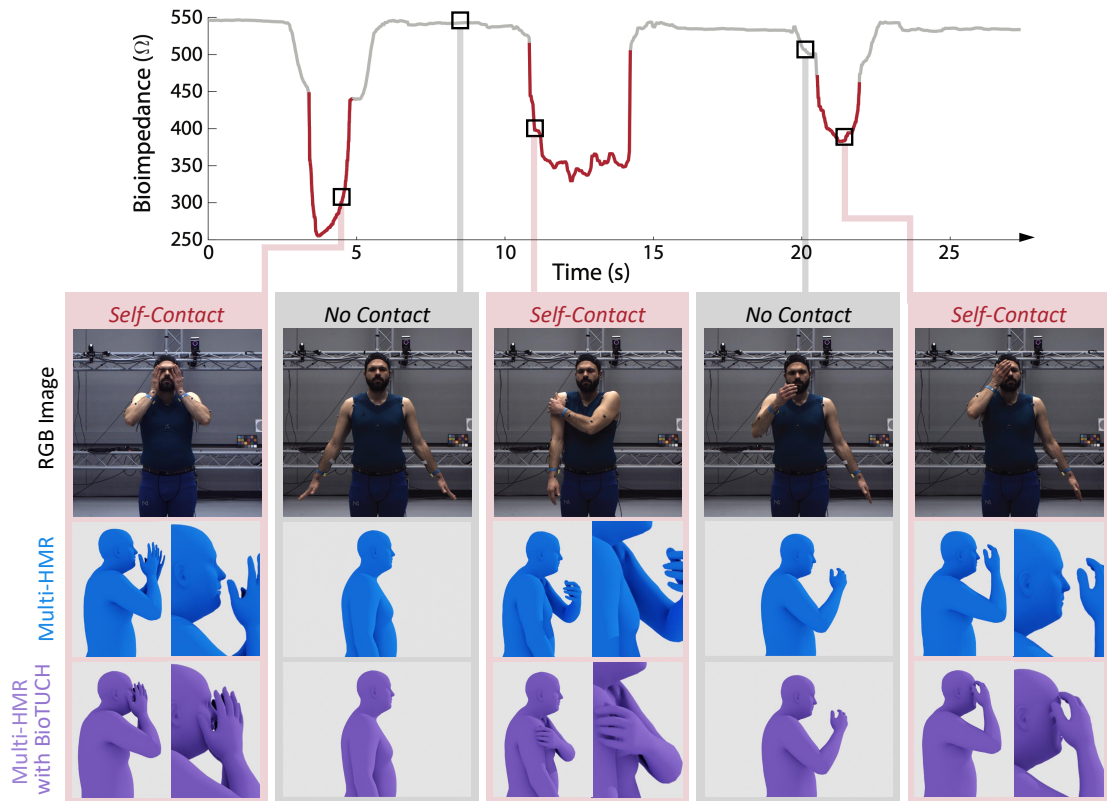


Figure 5.1: Human pose estimation methods struggle to reconstruct self-contact along the camera’s viewing axis. BioTUCH uses sharp changes in the bioimpedance signal measured between the wrists to estimate the beginning and end of such contacts. For all frames between these points (red segments of the signal), BioTUCH optimizes the results of off-the-shelf methods (such as Multi-HMR) to create plausible contact. Five selected poses (at the indicated time points) and their reconstructions are shown, with zoomed-in views of the contact regions.

accuracy of human-annotated contact labels. These factors influence the performance of the resulting regressor, TUCH [67]. Fieraru *et al.* [22] approach the data problem through manual annotation of two datasets. First, in the HumanSC3D dataset, subjects reproduce specific contact types (not full poses) while being recorded with four multi-view RGB cameras, and annotators then manually label the body regions in contact and their correspondence. Second, FlickrSC3D involves manual classification of web images into three categories: contact, no contact, and uncertain contact. The inclusion of an “uncertain” class highlights the inherent difficulty in obtaining definitive GT for self-contact poses from images alone. Multi-view capture setups could help, but their complexity and cost make them impractical for large-scale data collection. Finally, no prior work addresses the temporal dynamics of how self-contact states evolve over time.

The avatars resulting from standard HPE methods often have inaccurate arm poses, as the lack of hand contact propagates errors immediately to the wrist pose and consequently to the full arm pose. To address this cascading error problem, our approach allows refinement of the entire arm pose estimates when bioimpedance detects contact. BioTUCH optimizes more aggressively when visual ambiguity is highest: when contact occurs along the viewing direction. Consider a hand positioned in front of the face, as shown in Fig. 5.1; the proximity of the hand to the body is not directly observable, and the classical depth ambiguity means that contact is hard to resolve. We formulate a loss accounting for this uncertainty: when our sensor indicates contact but the 3D pose shows none, optimization drives the hand to contact the body while allowing greater freedom of movement along the  $z$ -axis, as it is most challenging to estimate from monocular images. We refine arm pose through selective joint optimization, updating only relevant joints (*i.e.*, shoulder, elbow, wrist) via a masked gradient strategy.

To evaluate our approach, we introduce a proof-of-concept dataset of synchronized RGB video, bioimpedance measurements, and 3D motion-capture (mocap) ground truth (GT) from three subjects. This in-lab dataset comprises 82 self-touch gestures derived from observational studies and neuropsychological research. We also include nine adversarial non-contact gestures to test detection robustness. Quantitative analysis demonstrates that our method achieves an average relative improvement of 11.67% in reconstruction accuracy during contact compared to visual pose estimators alone. To show applicability outside the lab, we present a miniature wearable bioimpedance sensor and perform an additional validation study; this sensor enables the capture of self-contact datasets with everyday clothing in natural indoor and outdoor settings.

By introducing the scalable and reliable framework of BioTUCH, we advance the field’s ability to generate contact-aware training data and more accurately reconstruct natural human movements involving self-contact. Specifically, our contributions include (1) a novel approach to sensing self-contact using bioimpedance; (2) a novel computer-vision method that integrates this contact information into the 3D human pose estimation task; (3) a new dataset that includes video with self-touch gestures, bioimpedance measurements, and GT 3D poses; and (4) the design of a miniature bioimpedance sensor that enables practical capture of self-contact in the wild. Our code and motion-capture dataset are available for research purposes at <https://biotuch.is.tue.mpg.de>

## 5.1 Method

BioTUCH is a novel framework that leverages bioimpedance sensing to optimize 3D arm-pose estimates during self-contact events. It focuses on dynamic contact detection during natural movement, using the optimal parameters established in Chapter 4.



Figure 5.2: Wearable setup. We created a miniaturized sensor that consists of a processor, shown here attached to the wearer’s shirt, and two electrodes embedded in the blue bracelets on their wrists. All components can be hidden beneath clothing.

### 5.1.1 Bioimpedance Sensing

As discussed in Chapter 4, bioimpedance sensing particularly excels at detecting skin-to-skin contacts, such as those shown in Fig. 5.1. While bioimpedance sensing can detect general self-contact between any body parts with appropriate electrode placement, this work focuses on self-touch due to its prevalence in SL. Following Chapter 4, we use the high-quality commercial impedance analyzer (MFIA, Zurich Instruments) with a wrist-to-wrist electrode configuration. However, this device’s dimensions ( $29\text{ cm} \times 23\text{ cm} \times 10\text{ cm}$ ), weight ( $\sim 4\text{ kg}$ ), and cost ( $\sim 15,000\text{ USD}$ ) make it impractical for large-scale data collection. Thus, we create a compact, low-cost, and wearable version [27] (see Fig. 5.2). Our custom device is small ( $2\text{ cm} \times 1.8\text{ cm} \times 1.1\text{ cm}$ ), lightweight ( $0.02\text{ kg}$ ), affordable ( $\sim 20\text{ USD}$ ), and easy to build and use. It operates on a small battery that allows over three hours of continuous use. The core component is a microcontroller (Adafruit QT Py ESP32-C3) mounted on a custom circuit board and connected via thin wires to two electrodes. To achieve secure and comfortable mounting on the wrists, we employ commercial electrode wristbands commonly used to prevent electrostatic discharge. This miniaturized device makes it possible to detect self-contact even when the subject is outdoors and wearing everyday clothing, demonstrating the feasibility of a more wearable and scalable solution.

### 5.1.2 BioTUCH

BioTUCH consists of two steps: (1) detecting self-contact using the bioimpedance signal and (2) optimizing the arm pose to enforce contact, if detected.

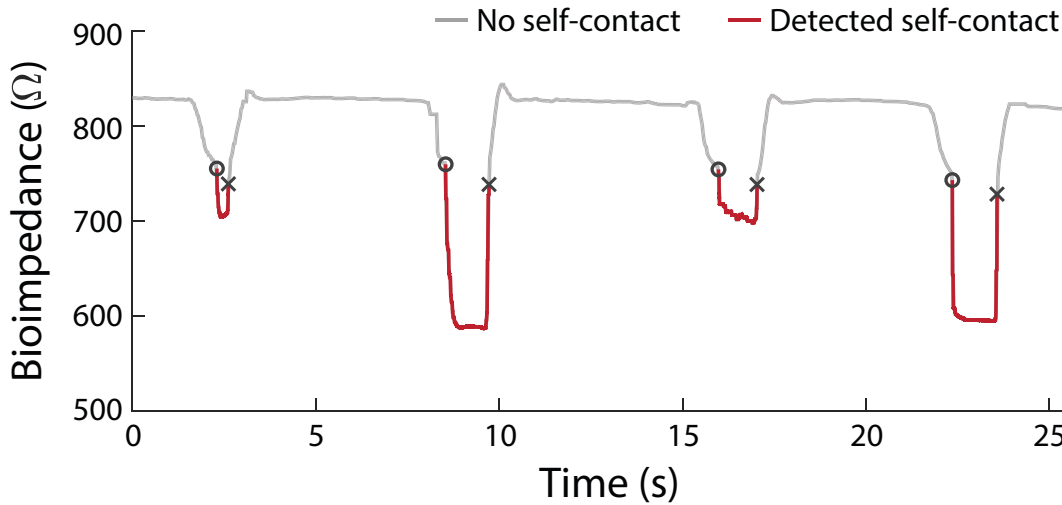


Figure 5.3: Sample bioimpedance signal showing four self-contact gestures over time. The bioimpedance magnitude begins to change when the person moves from the resting pose. The start of each self-contact triggers an abrupt drop in bioimpedance that is detected by our algorithm and marked with  $\circ$ . The thick red segments show the duration of the self-contacts, and their estimated ends are marked with  $\times$ .

### Self-Contact Detection

As shown in Fig. 5.3, wrist-to-wrist bioimpedance decreases sharply during instances of skin-to-skin contact. For self-contact detection, the signal is first resampled via linear interpolation to ensure a fixed sampling rate, then smoothed using a median filter with a window size of 100 ms, and finally differentiated. The beginning of each self-contact event is identified using an adaptive thresholding method applied to the processed signal, similar to techniques used in audio onset detection [6] and in detecting movement or muscle contraction [10, 38]. Specifically, we set the self-contact threshold to approximately one-third of the average of the three lowest minima in the processed recording; this thresholding procedure was designed and refined empirically. The end of each self-contact is defined as the time at which the bioimpedance magnitude returns to 98% of its pre-contact value. To reduce false positives, we apply temporal constraints informed by prior work on self-contact dynamics [64]. Our approach prioritizes specificity over sensitivity, allowing more false negatives (*i.e.*, missed contacts) than false positives (*i.e.*, incorrect contact detections) in order to reduce the risk of incorrectly optimizing non-contact poses. Finally, the detected contact events are converted into a binary touch sequence aligned with the video frames.

### Arm-Pose Optimization

We adopt the SMPL-X model [71] to represent the 3D body and output a 3D body mesh. We first run an off-the-shelf method (*e.g.*, Multi-HMR [2]) to obtain the initial SMPL-X estimates for a video sequence. We then average the estimated shape parameters ( $\beta$ ) and do not optimize body shape to maintain a consistent body shape across the sequence. In the first frame, we optimize body translation and global orientation (*i.e.*, pelvis joint) following SMPLify-X [71] and use the estimated values for the full sequence. We do not optimize the body pose for frames in which the bioimpedance-based sensor detects no contact. For frames with detected contact, BioTOUCH refines the arm poses through the following contact-aware optimization. For each hand, it identifies potential contact regions by computing distances between hand vertices  $\mathbf{v}$  and target vertices  $\mathbf{u}$  on the input SMPL-X mesh. Target vertices include all upper-body vertices (including the head, face, and opposite hand) apart from the hand’s own arm. Since existing methods perform reasonably well in the camera’s image plane ( $x$ - and  $y$ -coordinates) but struggle with depth ( $z$ -axis), distances from contact regions along the camera’s  $z$ -axis are weighted less during contact region detection, *i.e.*, a 1 cm distance along  $z$  is equivalent to a 0.25 cm distance in the viewing plane for determining potential contact pairs. This weighting is based on quantitative analysis that showed that errors along the camera’s viewing direction are typically 3–4 times larger. From these weighted distances, we determine the closest vertex pair(s):  $\mathcal{P}_h = \{(\mathbf{v}_i, \mathbf{u}_i)\}$ . When the hand is close to multiple body regions (*e.g.*, it could touch both the face and the opposite hand),  $\mathcal{P}_h$  contains multiple vertex pairs from different regions that are combined during optimization to accommodate multiple simultaneous contacts. The initial 3D positions of the vertices of  $\mathcal{P}_h$  are stored and used to maintain spatial consistency during optimization.

The weighted vertex-pair distances of both hands are then compared to determine which arms to optimize: both arms when the relative difference between distances is small (specifically, when  $\frac{d_{\max} - d_{\min}}{d_{\min}} \leq 0.5$ ), indicating both hands are similarly close to contact, and otherwise only the arm with the smaller distance. Given the selected arm(s), the relevant joint parameters (*i.e.*, shoulder, elbow, and wrist) are optimized through masked gradient updates:

$$\theta_{i+1} = \theta_i - \eta \nabla_{\theta} \mathcal{L} \odot \mathbf{M}_a \quad (5.1)$$

where  $\theta_i \in \mathbb{R}^{63}$  represents the SMPL-X body pose parameters in axis-angle representation at iteration  $i$ ,  $\eta$  is the learning rate,  $\nabla_{\theta} \mathcal{L}$  is the gradient of the loss function with respect to all pose parameters, and  $\mathbf{M}_a \in \{0, 1\}^{63}$  is the binary mask. The loss function through which we update the arm pose is:

$$\mathcal{L} = \mathcal{L}_{2D} + \lambda_{\text{contact}} \mathcal{L}_{\text{contact}} \quad (5.2)$$

where  $\mathcal{L}_{2D}$  ensures consistency with the 2D joint observations of the arms, and

$$\mathcal{L}_{\text{contact}} = \mathcal{L}_{\text{consistency}} + \mathcal{L}_{\text{interpenetration}} + \mathcal{L}_{\text{proximity}} \quad (5.3)$$

$\mathcal{L}_{\text{consistency}}$  preserves the initial spatial relationships,  $\mathcal{L}_{\text{interpenetration}}$  prevents mesh interpenetration through a barrier function, and  $\mathcal{L}_{\text{proximity}}$  encourages contact convergence. Specifically,

$$\mathcal{L}_{\text{proximity}} = \sum_{h \in \mathcal{H}} \sum_{(\mathbf{v}_i, \mathbf{u}_i) \in \mathcal{P}_h} \sum_{d \in \{x, y, z\}} \omega_d |v_i^d - u_i^d| \quad (5.4)$$

where  $\mathcal{H}$  is the set of active hands and  $\omega_x$ ,  $\omega_y$ , and  $\omega_z$  are camera-adaptive weights (with depth weighted four times higher than the viewing plane to aggressively minimize contact errors along the camera’s viewing direction). Optimization continues until contact is achieved (distances  $\leq 5$  mm in all axes) or iteration limits are reached. This approach achieves physically plausible self-contact while maintaining consistency with both 2D keypoint observations and the initial pose estimates. Finally, for temporal consistency, we apply OneEuroFilter-based smoothing in post-processing, as in [45] (see supplementary video at <https://biotuch.is.tue.mpg.de>).

## 5.2 Dataset

To evaluate BioTUCH, we collected a new dataset of synchronized frontal RGB videos, bioimpedance measurements, and 3D motion capture. Figure 5.4 shows five sample poses from this dataset, representing the first collection of directly sensed contact-aware training data suitable for improving pose estimation models. The experimental procedure was reviewed and approved by our local ethics board.



Figure 5.4: Sample contacts that cause a significant change in the bioimpedance of the user. Any self-contact (direct on skin or through body hair) greatly impacts the wrist-to-wrist bioimpedance signal. The videos from which these images were taken and their ground-truth SMPL-X meshes are part of our newly collected dataset.

The dataset comprises 82 *common dynamic self-touch gestures* derived from observational studies, neuropsychological research, and poses from the MTP dataset [67]. These gestures are organized into the following categories according to the body parts in contact:

1. 38 hand-to-face or hand-to-head poses (*e.g.*, temple touch, nose touch, hands covering the face) reflecting attentional refocusing [64] or emotional regulation [34, 70];
2. 12 hand-to-hand gestures (*e.g.*, hand wringing, hands clasped, fingers interlocked) associated with emotional regulation [67, 70] and cognitive load [67];
3. 11 hand-to-neck or hand-to-shoulder movements (*e.g.*, shoulder rub) linked to self-support and stress relief [34];
4. Nine hand-to-chest or hand-to-torso movements (*e.g.*, palm over heart, self-hugs) often associated with self-reassurance or stress relief [16];
5. Eight behind-the-back gestures (*e.g.*, holding hands or clasping the opposite forearm behind the back) related to postural control, tension release, or self-soothing (note that these contacts are not visible in the frontal RGB camera);
6. Four hands-to-body gestures, where both hands contact different body parts simultaneously (*e.g.*, one hand supporting the head and the other touching the other arm), typical of postural support.

We also include nine adversarial non-contact gestures that can be visually interpreted as self-touch contacts (*e.g.*, hand close to the forehead or mouth) when seen in a frontal view. For unilateral gestures, we include variants performed with each hand.

Three participants (two females and one male) took part in the study. Although the number of participants is limited, our electrode configuration and bioimpedance sensing parameters are based on Chapter 4, which validated this sensing method across a diverse population that varied in body mass index, sex, ethnicity, and handedness. Based on our prior results, we measured the bioimpedance magnitude using an alternating current signal at 2.29 MHz, with data sampled at 13.4 kHz. The participants were captured with a Vicon mocap system at 30 fps while wearing the two electrode bracelets. The mocap system was synchronized with the impedance analyzer (MFIA, Zurich Instruments) and a frontal 1300×1400 RGB camera at 30 fps, framing a full-body view. All participants performed the gestures while standing; two participants repeated them while sitting. The arms rested at the participant’s sides at the beginning and end of each gesture. GT SMPL-X meshes were obtained by scanning the participants in a 4D body scanner in several poses and fitting the SMPL-X model to these scans. A personalized body-shape mesh of each participant was obtained by averaging the SMPL-X meshes in the canonical pose. MoSh++ was then used to fit this mesh to the mocap markers [57],

Dataset	# of Contact Frames	Data Type	Motion	Annotation
3DCP [67]	1,653	SMPL-X (from 3D scans, 3D mocap)	Static	Automatic (geometry)
MTP [67]	3,731	RGB images, (fitted) SMPL-X	Static	Automatic (geometry)
DSC [67]	30,000	RGB images	Static	Manual
HumanSC3D [22]	4,128	3D mocap, multi-view RGB images	Dynamic	Manual
FlickrSC3D [22]	3,969	RGB images	Static	Manual
Ours	19,183	SMPL-X (from 3D mocap), RGB images, bioimpedance	Dynamic	Automatic (bioimpedance)

Table 5.1: Comparison of datasets for self-contact.

providing an accurate GT 3D body mesh for each frame (Fig. 5.4). Table 5.1 compares our dataset with existing datasets related to self-contact.

## 5.3 Experiments

### 5.3.1 Evaluation of Self-Contact Detection

To quantitatively evaluate the contact-detection accuracy of our sensing approach, we manually labeled binary self-contact over time for 57 gestures (including adversarial non-contact gestures); these were randomly selected from all three participants and both sitting and standing poses, excluding gestures used for setting the algorithm’s parameters. It is important to note that both frames and bioimpedance were used for manual labeling, as frames alone were not reliable in several cases. The algorithm has a sensitivity of 0.858 and specificity of 0.992 (false-negative rate of 0.142 and false-positive rate of 0.008). These high sensitivity and specificity values confirm the reliability of the sensor as a proxy for GT in contact detection. Furthermore, none of the non-contact gestures were misclassified as contact.

### 5.3.2 Quantitative Evaluation

We quantitatively evaluate BioTUCH, comparing vanilla Multi-HMR [2], AiOS [84], and TUCH [67] (converted to SMPL-X) with their BioTUCH extensions. Multi-HMR and AiOS are the two state-of-the-art methods in human pose estimation, and TUCH

Method	PA-V2V ↓ (mm)	Euclidean Joint Location Error ↓ (mm)			Detection Rate ↑ (%)	V-Distance ↓ (mm)
		Shoulder	Elbow	Wrist		
Multi-HMR [2]	57.46	<b>23.49</b>	<b>29.86</b>	65.37	41.28	87.48
+ 2D Loss	77.41	26.74	39.31	84.97	42.16	124.66
+ Contact Loss	50.74	23.71	31.71	57.09	<b>79.35</b>	72.59
+ BioTUCH	<b>50.21</b>	23.95	30.99	<b>56.29</b>	78.34	<b>71.22</b>
AiOS [84]	72.24	<b>21.84</b>	<b>39.83</b>	77.29	45.87	99.02
+ 2D Loss	86.61	27.01	41.31	93.38	41.58	130.64
+ Contact Loss	65.10	22.13	43.99	68.29	<b>80.60</b>	82.01
+ BioTUCH	<b>62.79</b>	22.42	40.65	<b>65.61</b>	78.48	<b>79.16</b>
TUCH [67]	70.55	<b>28.24</b>	41.03	58.71	59.46	96.31
+ 2D Loss	77.12	33.66	49.00	77.03	49.52	110.46
+ Contact Loss	64.95	28.24	40.65	53.24	<b>85.25</b>	86.95
+ BioTUCH	<b>63.99</b>	28.27	<b>40.59</b>	<b>52.91</b>	84.60	<b>86.26</b>

Table 5.2: Performance evaluation reporting vertex-to-vertex arm mesh error (PA-V2V), Euclidean error for each arm joint, and the contact metrics of detection rate and distance between the contact vertices in GT and in the estimate. We compare each off-the-shelf method with its BioTUCH extension, including variants that isolate the contributions of BioTUCH’s two losses: 2D loss and contact loss. For each input method, the best result among its variants is shown in bold font. ↓ lower value means better performance. ↑ higher value means better performance.

is tuned for self-contact. The evaluation uses synchronized meshes from Vicon markers [57] as GT and is conducted on a per-frame basis before any temporal smoothing. All methods estimate SMPL-X meshes with the same topology, enabling Procrustes alignment (PA) with the GT for fair comparison. We compute the mean per-vertex error (PA-V2V) of the optimized full arm regions and the Euclidean errors of the key arm joints, as well as the contact accuracy. Specifically, we compare each method’s output with the GT meshes to compute detection rate (*i.e.*, how often a GT contact is reconstructed) and contact preservation (*i.e.*, how well spatial relationships between contacting vertices are maintained) by measuring the Euclidean distances between the same vertex pairs that were in contact in the GT. Table 5.2 reports the results of these metrics for the three input methods, including an evaluation showing how each loss component (*i.e.*, 2D loss and contact loss) contributes to performance.

Among the input methods, AiOS has the highest PA-V2V error (72.24 mm), while Multi-HMR achieves the best baseline performance (57.46 mm). As expected, TUCH has the highest contact detection rate (59.46%) among the baselines, reflecting its specialization for self-contact scenarios. In general, BioTUCH consistently improves all input methods across key metrics. For PA-V2V error, we achieve an average 11.67% improvement. Joint-level analysis reveals varied improvements across the arm: shoulder

and elbow joint errors show overall a slight degradation (0.36 mm and 0.50 mm), while wrist joint locations drastically improve (8.85 mm). BioTUCH increases contact detection rates considerably across all methods, with an average increase of 31.60 percentage points and an average reduction of 15.39 mm in distances between contacting vertex pairs. These findings indicate that our method not only reconstructs contacts more reliably but also maintains better spatial relationships between contacting surfaces.

The analysis of the loss components reveals that while combining 2D and contact losses improves baseline performance, the contact loss provides the dominant contribution. The 2D loss (SMPLify-X’s re-projection error) shows minimal benefit compared to our contact loss component. The highest detection rate is achieved when adding the contact loss alone, as it is not constrained by keypoint positions; however, spatial accuracy of the contact decreases compared to the full BioTUCH loss, as without the 2D loss the optimization prioritizes achieving contact and does not preserve alignment with the original keypoints. This result confirms that bioimpedance provides unique information that cannot be derived from 2D keypoints alone, validating our multi-modal approach of prioritizing contact-aware optimization. Consistent with its reliable baseline performance, Multi-HMR achieves the best overall results when enhanced with BioTUCH, reaching 50.21 mm PA-V2V error.

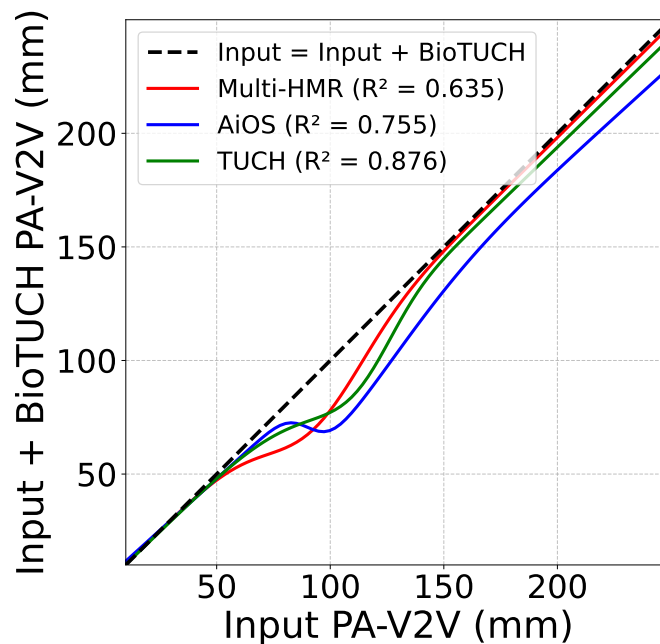


Figure 5.5: Robustness evaluation. The plot shows PA-V2V errors before ( $x$ -axis) and after ( $y$ -axis) applying BioTUCH to the three input methods. Points below the dashed black line indicate improvement. The  $R^2$  values show how well the fitted curves explain the variance in the data.

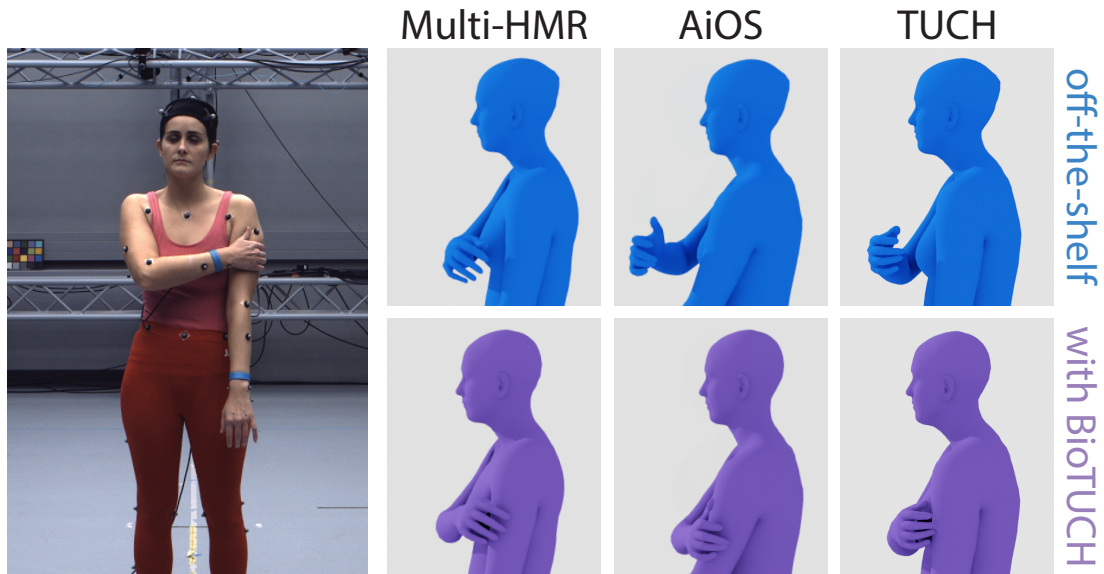


Figure 5.6: Qualitative results. The left column shows an RGB image from a sample self-contact gesture in our dataset. On the right, the estimates of the off-the-shelf methods (*i.e.*, Multi-HMR, AiOS, TUCH) are shown in blue in the first row, and their BioTUCH extensions are shown in purple in the second row. BioTUCH clearly optimizes the arm joints to bring the hand into self-contact at the observed location in the image plane.

Finally, we conduct a robustness evaluation to investigate BioTUCH’s effectiveness across different levels of input reconstruction error. Figure 5.5 shows the relationship between the PA-V2V errors of the input methods and the corresponding errors after applying BioTUCH. The analysis reveals that BioTUCH consistently improves all baseline methods across the entire error spectrum, with the most pronounced improvements in the 50–150 mm range.

### 5.3.3 Qualitative Evaluation

We qualitatively evaluated the three input methods and their BioTUCH extensions on our dataset. Among the input methods, we noticed that TUCH performs particularly well for hand-to-face contacts. Multi-HMR and AiOS instead do not exhibit any significant variations across different body parts. Overall, as expected, the most accurate contact predictions for all input methods occur when the contact is along the camera’s  $xy$  plane (*e.g.*, a hand touching the temple), and the most challenging scenarios are the behind-the-back gestures. In both cases, BioTUCH has minimal impact. In case of contacts in the  $xy$  plane, the contact is often already reconstructed by the input methods, so BioTUCH does not optimize the 3D pose. Conversely, the significant errors of the estimated joint positions of the behind-the-back gesture do not allow meaningful con-



Figure 5.7: Results of two in-the-wild captures. When Multi-HMR fails to reconstruct a self-contact detected by our miniature bioimpedance sensor, BioTUCH optimizes the arm joints to resolve depth ambiguity and enforce contact.

tacts. BioTUCH proves particularly beneficial for visible contacts that occur in line with the camera’s optical axis. Although these contacts are visible, they pose challenges for the off-the-shelf methods due to the depth ambiguity problem. Figure 5.6 illustrates the estimates for a sample gesture from Multi-HMR, AiOS, and TUCH, along with the improvements achieved by BioTUCH. Overall, BioTUCH effectively optimizes arm poses to ensure contact while maintaining the consistency of 2D points. However, the optimization halts once contact is achieved. Consequently, inaccuracies in the finger articulation of the input method can affect the stopping condition of BioTUCH, as seen for TUCH in Figure 5.6. More qualitative results are shown in the supplemental video at <https://biotuch.is.tue.mpg.de>

### 5.3.4 In-the-Wild Capture

To demonstrate outside-the-lab feasibility, we collected two short in-the-wild sequences using our miniature sensor (see Fig. 5.2) worn beneath clothing and a single frontal or side-view camera. The bioimpedance data and the video were recorded and automatically synchronized using MATLAB. As shown in Fig. 5.7, Multi-HMR fails to accurately reconstruct self-contact along the camera axis in both cases. Whether the camera is viewing the front or side of the user, BioTUCH addresses the depth ambiguity by optimizing

the arm pose, successfully generating the detected contact.

### 5.3.5 Application to Sign Language

For SLC, we use SGNify’s estimates as input to BioTUCH. SLC poses particular challenges for contact detection due to the frequent occurrence of self-touch, where many signs involve contact between the hands or between the hands and the body (for more details, please refer to Sec. 2.1.3). As demonstrated in Fig. 5.8, BioTUCH refines the arm poses to establish contact where none was originally detected by SGNify. As such, BioTUCH addresses SGNify’s remaining challenge with contact-dependent signs, which SGNify and standard HPE methods struggle to capture accurately due to depth ambiguity in monocular video and the insufficient contact-specific training data.

## 5.4 Discussion

The success of BioTUCH builds directly on our findings from the previous chapter, demonstrating that wrist-to-wrist bioimpedance measurements can indeed provide reliable pGT contact information for poses involving self-touch. Our quantitative experiments demonstrate that BioTUCH improves contact reconstruction in all tested off-the-shelf methods by leveraging direct-contact sensing. Our qualitative evaluation further demonstrates that BioTUCH produces visually improved reconstructions for general HPE and SLC, with more accurate hand-to-hand and hand-to-body contact poses, overcoming the limitations observed in SGNify’s results. Extending BioTUCH to detect contact between additional body parts through strategic electrode placement would broaden

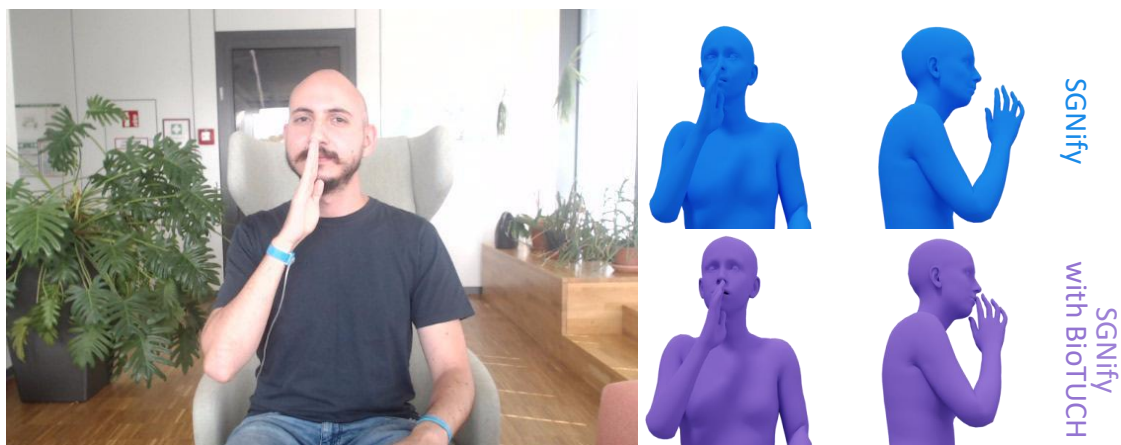


Figure 5.8: The sign  $\mu\epsilon\sigma\eta\mu\epsilon\rho\iota$  in Cypriot Sign Language reconstructed by SGNify and SGNify with BioTUCH. The side view confirms that BioTUCH correctly establishes hand-to-face contact.

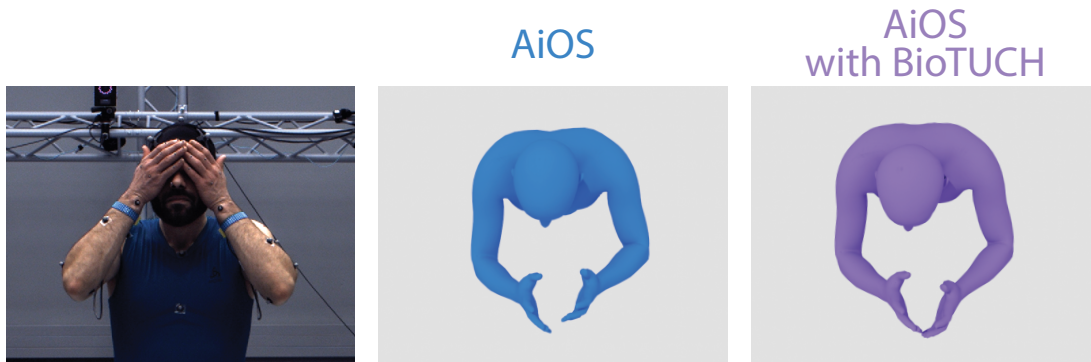


Figure 5.9: Input RGB image for a challenging two-handed self-touch gesture, with top-down views of the estimates of AiOS and AiOS with BioTUCH.

its applicability beyond self-touch and SLC, enabling new use cases in healthcare monitoring, behavioral research, and human-computer interaction systems that require precise contact detection.

#### 5.4.1 Limitations and Future Work

Several areas present opportunities for future enhancement. Our current approach relies on the input mesh estimates to determine which parts of the body should be in contact. However, as shown in Fig. 5.9, input meshes can exhibit substantial errors along all axes, leading to suboptimal contact region identification. Additionally, when a single hand contacts multiple body parts simultaneously (as in Fig. 5.9, where each hand touches both the other hand and the face), distance-based region identification is fundamentally limited by the initial mesh’s spatial configuration. For instance, in Fig. 5.9, errors along the  $z$ -axis are significantly larger than those in the  $xy$ -plane, and BioTUCH prioritizes optimizing the hand-to-hand contact over the hand-to-face contacts. Even localized inaccuracies, such as incorrect finger articulation, can stop the optimization too early, as observed in the bottom row of Fig. 5.7. Future work should address these limitations by using visual evidence to identify the contact regions when contact is visible in the raw input images, thus allowing the optimization to maintain all visually identified contacts rather than selecting only the closest mesh-based estimate. Vision-Language Models (VLMs) present a promising approach for this task. By leveraging their pre-trained understanding of human anatomy and spatial relationships, they can determine the contacting body parts directly from RGB images when bioimpedance detects contact.

Beyond these improvements, BioTUCH can be readily extended in several directions. We focused on hand-initiated contacts as they are most common, but bioimpedance sensing can be used to detect self-contact at any location by repositioning the electrodes (*e.g.*, moving them to the ankles for lower-limb contact detection) and updating the active

body parts and target vertices. Furthermore, we used binary contact detection; however, bioimpedance signals can provide richer information, including the location and size of the contact (see Chapter 4). This sensing approach can also detect skin-to-clothing contact through fabric when both fabric sides contact the skin (see Chapter 4), as we observed in some hand-to-chest contacts in our dataset. Finally, large-scale collection with our miniature sensor could provide improved training data for methods that regress body parameters from RGB images. The next step would be to collect such a dataset and train a regressor on bioimpedance-corrected meshes, similar to how TOUCH [67] trained on visually derived pGT. These directions, along with broader reflections on this dissertation’s contributions, implications, and limitations, are discussed in the concluding chapter that follows.

# Chapter 6

## Discussion

This dissertation addressed a fundamental challenge in developing accessible technology for the Deaf community: the lack of methods that enable accurate SLC at scale. While existing methods achieve reasonable performance on typical human motion, they struggle with the complex spatial-temporal dynamics of SL, often resulting in avatars that have poses that are physically plausible but linguistically incorrect. Addressing the limitations of current methods is crucial not only for ensuring equitable access to technology for the millions of Deaf people worldwide who use SL as their primary mode of communication but also for advancing motion capture technology.

### 6.1 Summary of Contributions

This dissertation significantly advanced SLC through three interconnected contributions that integrate insights from linguistics and bioimpedance sensing into vision-based HPE. The first contribution established that phonological constraints governing sign formation can be encoded as differentiable loss functions within HPE. SGNify formalizes universal linguistic rules, *i.e.*, hand-pose symmetry and hand-pose invariance, which are derived from phonological conditions, into mathematical constraints that improve hand pose estimation accuracy. The generalization of this approach to different SLs enables its application to SLs with scarce resources, where it would be impossible to have enough data to train an SL-specific regressor. The perceptual evaluation revealed that linguistic constraints alone could bridge the gap between automated reconstruction and human-level sign intelligibility, achieving recognition rates (86.2%) statistically indistinguishable from source videos (90.9%). This level of performance enables automated systems to produce signing avatars of sufficient quality to support accessibility applications.

The second contribution addressed a fundamental limitation of monocular HPE methods: reconstructing self-touch. Through systematic investigation of 27 self-touch poses across 30 diverse participants, we established that wrist-to-wrist bioimpedance magnitude measured at high frequencies (237.8 kHz to 4.1 MHz) reliably distinguishes skin-to-skin self-touch from no-contact poses. All skin-to-skin poses were detected with  $p < 0.001$  statistical significance, confirming contact specificity across diverse individuals regardless of sex, ethnicity, or BMI. This systematic characterization established

the fundamental sensing parameters necessary for practical implementation and enabled development of a miniaturized and cheap sensor (2 cm × 1.8 cm × 1.1 cm, 0.1 kg, ~20 EUR) that demonstrates feasibility for real-world deployment.

The third contribution showed how temporal contact information from bioimpedance can be systematically integrated with spatial localization from computer vision to create physically plausible contact; knowing *when* contact occurs resolves visually ambiguous scenarios. The optimization strategy selectively refines arm poses only during frames where contact is detected. Quantitative evaluation across multiple human pose estimation methods (*i.e.*, Multi-HMR, AiOS, TOUCH) demonstrated an average relative improvement of 11.67% in reconstruction accuracy with an absolute improvement of 8.85 mm in wrist joint localization, *i.e.*, the terminal point of the kinematic chain needed for self-touch, and an absolute increase of 31.60 percentage points in correctly detecting contact when it is present in the ground truth. Since BioTOUCH functions as a post-processing refinement that enhances any HPE method, it can be straightforwardly integrated with future methods.

## 6.2 Implications and Applications

The demonstrated technical capabilities enable applications that directly support compliance with accessibility legislation such as the Americans with Disabilities Act [87] and the European Accessibility Act [19], while establishing methodological principles for broader human motion analysis. SL accessibility is supported through multiple pathways. In educational applications, accurate 3D signing avatars could enable learners to examine signs from multiple viewpoints, supporting spatial understanding of complex hand poses, movements, and self-touch, an advantage that 2D representations cannot provide. The 3D format also enables view-independent rendering for virtual and augmented reality applications. Moreover, it is ideal for capturing the spatial nature of co-articulation between signs. Phonologically precise 3D motions with natural co-articulation, captured at scale using the demonstrated cost-effective approach, provide the data needed to train high-quality SL generators.

Beyond these direct applications, this dissertation establishes several methodological advances for HPE research that extend beyond SLC. The demonstrated success of linguistic constraint formalization suggests applicability to other rule-governed movement systems. For example, dance notation systems could similarly guide reconstruction algorithms, where choreographic principles would constrain pose sequences in the same way as phonological rules constrain sign formation. Furthermore, bioimpedance-based contact detection can provide reliable pGT contact data at scale. Because bioimpedance sensing easily extends to other body regions through strategic electrode positioning, BioTOUCH's optimization approach can be straightforwardly adapted to diverse contact scenarios. The resulting pGT data can train VLMs to more accurately detect and interpret contact in human movement, enabling researchers in healthcare monitoring (*e.g.*, de-

tecting patient self-contact behaviors), sports analysis (*e.g.*, tracking athlete self-contact patterns), and extended reality systems (*e.g.*, improving avatar realism) to investigate contact-dependent phenomena that were previously difficult to study systematically.

## 6.3 Limitations and Future Work

To support practical deployment and maximize accessibility, future work should address the current limitations across technical, methodological, and community dimensions. From a technical perspective, both SGNify and BioTUCH are optimization-based approaches, which, while effective for offline reconstruction and dataset generation, are computationally intensive and unsuitable for real-time applications. A transition toward regression-based architectures trained on large-scale datasets is a promising direction, made feasible by the cost-effective and scalable sensing pipeline developed in this work. This framework enables data collection at the scale necessary for robust training. However, before practical deployment, constraints related to bioimpedance sensing robustness must be addressed. While the miniaturized bioimpedance sensor performed robustly under controlled conditions, its sensitivity to environmental factors such as temperature and humidity has yet to be evaluated. Moreover, the sensor's performance under varying clothing conditions was only preliminarily analyzed, limited to coarse material categories in our 30-participant study, and requires investigation under the varied conditions users encounter in natural environments. Furthermore, ergonomic refinements are needed to ensure the sensor does not interfere with the naturalness of signing, given that small constraints on movement can affect sign articulation.

In this dissertation, we reported all relevant metrics across disciplines: HPE validation emphasizes geometric accuracy, measured using the Vicon mocap system; linguistic validation prioritizes sign recognition and naturalness, assessed through user studies with Deaf participants; and sensor validation focuses on detection reliability, determined via manual labeling. Beyond the substantial effort required to compute each of these metrics, our findings also revealed that evaluation frameworks from different disciplines often conflict. For example, SGNify's phonological constraints occasionally produced signs that had higher V2V error compared to other methods, yet were more intelligible to Deaf users in linguistic evaluations. These inconsistencies underscore the need for integrated evaluation frameworks that reflect the multifaceted nature of SLC.

Finally, although this research involved Deaf community members during dataset development and avatar evaluation, their participation primarily occurred during validation stages. Feedback from these interactions strongly indicated the importance of involving the Deaf community earlier and more deeply in the development process. A co-design approach, where Deaf people are actively involved in shaping system goals, constraints, and evaluation criteria, would better ensure that the resulting technology aligns with linguistic norms, cultural values, and community priorities, advancing inclusive and community-centered SLC.

## 6.4 Conclusion

This dissertation enables one of the crucial steps for inclusive access to information for the global Deaf community: the collection of high-quality 3D data at scale. Current HPE systems struggle with the complex visual-spatial modality of SL, creating technical and usability barriers to accessibility technology development. Through interdisciplinary integration, this research shows that these barriers can be overcome through technical innovation validated with comprehensive user studies. Three interconnected contributions validate this approach: SGNify reconstructs signs that proficient signers recognize at rates statistically indistinguishable from source videos (86.2% vs 90.9%) by formalizing linguistic domain knowledge as computational constraints; bioimpedance sensing provides reliable contact detection (85.8% sensitivity, 99.2% specificity) while achieving practical miniaturization to 20-EUR wearable devices, democratizing access to contact-detection technology; and BioTUCH produces substantial improvements to HPE (11.67% average) by resolving fundamental depth ambiguities through multimodal sensing. The path forward requires overcoming the challenges of practical deployment and continued collaboration between technologists and SL communities.

Returning to the opening example of this dissertation, the ASL sign THANKS, with its complex spatial-temporal dynamics and contact requirements, can now be effectively reconstructed, as shown in Fig. 6.1.



Figure 6.1: THANKS (ASL) in the original video and reconstructed by SGNify with BioTUCH. Both frontal and side views demonstrate effective reconstruction, with the hand contacting the face.

# Bibliography

- [1] Anguelov, D., Srinivasan, P., Koller, D., Thrun, S., Rodgers, J., and Davis, J. (2005). Scape: shape completion and animation of people. In *ACM Siggraph 2005 Papers*, pages 408–416.
- [2] Baradel, F., Armando, M., Galaaoui, S., Brégier, R., Weinzaepfel, P., Rogez, G., and Lucas, T. (2025). Multi-HMR: Multi-person whole-body human mesh recovery in a single shot. In *European Conference on Computer Vision*, pages 202–218. Springer.
- [3] Bartels, E. M., Sørensen, E. R., and Harrison, A. P. (2015). Multi-frequency bioimpedance in human muscle assessment. *Physiological Reports*, **3**(4), e12354.
- [4] Battison, R. (1978). *Lexical Borrowing in American Sign Language*. Education Resources Information Center (ERIC) .
- [5] Bilal, S., Akmeliawati, R., El Salami, M. J., and Shafie, A. A. (2011). Vision-based hand posture detection and recognition for sign language – a study. In *International Conference on Mechatronics (ICOM)*, pages 1–6.
- [6] Böck, S. and Widmer, G. (2013). Maximum filter vibrato suppression for onset detection. In *Proceedings of the International Conference on Digital Audio Effects (DAFx)*, volume 7, page 4. Citeseer.
- [7] Buendía, R., Gil-Pita, R., and Seoane, F. (2011). Cole parameter estimation from total right side electrical bioimpedance spectroscopy measurements—influence of the number of frequencies and the upper limit. In *Proceedings of the Annual International Conference of the IEEE Engineering in Medicine and Biology Society (EMBC)*, pages 1843–1846.
- [8] Campa, F., Toselli, S., Mazzilli, M., Gobbo, L. A., and Coratella, G. (2021). Assessment of body composition in athletes: A narrative review of available methods with special reference to quantitative and qualitative bioimpedance analysis. *Nutrients*, **13**(5), 1620.
- [9] Cao, Z., Simon, T., Wei, S.-E., and Sheikh, Y. (2017). Realtime multi-person 2D pose estimation using part affinity fields. In *Computer Vision and Pattern Recognition (CVPR)*, pages 1302–1310.

- [10] Carvalho, C. R., Fernández, J. M., Del-Ama, A. J., Oliveira Barroso, F., and Moreno, J. C. (2023). Review of electromyography onset detection methods for real-time control of robotic exoskeletons. *Journal of Neuroengineering and Rehabilitation*, **20**(1), 141.
- [11] Daněček, R., Black, M. J., and Bolkart, T. (2022). EMOCA: Emotion driven monocular face capture and animation. In *Computer Vision and Pattern Recognition (CVPR)*, pages 20311–20322.
- [12] de Talhouet, H. and Webster, J. G. (1996). The origin of skin-stretch-caused motion artifacts under electrodes. *Physiological Measurement*, **17**(2), 81.
- [13] Deurenberg, P., Weststrate, J. A., Paymans, I., and Van der Kooy, K. (1988). Factors affecting bioelectrical impedance measurements in humans. *European Journal of Clinical Nutrition*, **42**(12), 1017–1022.
- [14] Dezfuli, N., Khalilbeigi, M., Huber, J., Müller, F., and Mühlhäuser, M. (2012). PalmRC: Imaginary palm-based remote control for eyes-free television interaction. In *Proceedings of the European Conference on Interactive TV and Video*, pages 27–34, Berlin, Germany. ACM.
- [15] Dheman, K., Mayer, P., Eggimann, M., Schuerle, S., and Magno, M. (2021). ImpediSense: A long lasting wireless wearable bio-impedance sensor node. *Sustainable Computing: Informatics and Systems*, **30**, 100556.
- [16] Dreisoerner, A., Junker, N. M., Schlotz, W., Heimrich, J., Bloemeke, S., Ditzen, B., and van Dick, R. (2021). Self-soothing touch and being hugged reduce cortisol responses to stress: A randomized controlled trial on stress, physical touch, and social identity. *Comprehensive Psychoneuroendocrinology*, **8**, 100091.
- [17] Duarte, A., Palaskar, S., Ventura, L., Ghadiyaram, D., DeHaan, K., Metze, F., Torres, J., and Giro-i Nieto, X. (2021). How2Sign: A large-scale multimodal dataset for continuous american sign language. In *Computer Vision and Pattern Recognition (CVPR)*, pages 2735–2744.
- [18] Eccarius, P. and Brentari, D. (2007). Symmetry and dominance: A cross-linguistic study of signs and classifier constructions. *Lingua*, **117**(7), 1169–1201.
- [19] European Union (2019). Directive (EU) 2019/882 of the European Parliament and of the Council of 17 April 2019 on the accessibility requirements for products and services (European Accessibility Act). <https://eur-lex.europa.eu/legal-content/EN/TXT/?uri=CELEX%3A32019L0882>. Official Journal of the European Union L 151, 7 June 2019.

- [20] Evans, W. D., McClagish, H., and Trudgett, C. (1998). Factors affecting the in vivo precision of bioelectrical impedance analysis. *Applied Radiation and Isotopes*, **49**(5-6), 485–487.
- [21] Feng, Y., Choutas, V., Bolkart, T., Tzionas, D., and Black, M. (2021). Collaborative regression of expressive bodies using moderation. In *International Conference on 3D Vision (3DV)*, pages 792–804.
- [22] Fieraru, M., Zanfir, M., Oneata, E., Popa, A.-I., Olaru, V., and Sminchisescu, C. (2021). Learning complex 3D human self-contact. In *Proceedings of the AAAI Conference on Artificial Intelligence (AAAI)*, pages 1343–1351.
- [23] Filntisis, P. P., Retsinas, G., Paraperas-Papantoniou, F., Katsamanis, A., Roussos, A., and Maragos, P. (2023). SPECTRE: Visual speech-informed perceptual 3D facial expression reconstruction from videos. In *Computer Vision and Pattern Recognition Workshops (CVPRw)*, pages 5745–5755.
- [24] Forte, M.-P., Kulits, P., Huang, C.-H. P., Choutas, V., Tzionas, D., Kuchenbecker, K. J., and Black, M. J. (2023). Reconstructing signing avatars from video using linguistic priors. In *Computer Vision and Pattern Recognition (CVPR)*, pages 12791–12801.
- [25] Forte, M.-P., Athanasiou, N., Ballardini, G., Bartels, U. J., Kuchenbecker, K. J., and Black, M. J. (2025a). Contact-aware refinement of human pose pseudo-ground truth via bioimpedance sensing. In *International Conference on Computer Vision (ICCV)*, pages 5071–5080.
- [26] Forte, M.-P., Vardar, Y., Javot, B., and Kuchenbecker, K. J. (2025b). Dataset and code for “wrist-to-wrist bioimpedance can reliably detect self-touch”. Edmond, V1. [Online]. Available: <https://doi.org/10.17617/3.P0NEOF>.
- [27] Forte, M.-P., Kuchenbecker, K. J., Bartels, U. J., and Ballardini, G. (2025c). System und verfahren zur detektion von kontakt von teilen eines körpers. German Patent Application DE 10 2025 106 025.8 (provisional). Filed with the German Patent and Trademark Office on February 18, 2025, by Max-Planck-Gesellschaft zur Förderung der Wissenschaften e.V.
- [28] Forte, M.-P., Vardar, Y., Javot, B., and Kuchenbecker, K. J. (2025d). Wrist-to-wrist bioimpedance can reliably detect self-touch. In *IEEE Transactions on Instrumentation and Measurement*.
- [29] Geman, S. (1987). Statistical methods for tomographic image restoration. *Proceedings of the 46th Session of the International Statistical Institute, Bulletin of the ISI*, **52**, 5–21.

- [30] Glickman, M. E., Rao, S. R., and Schultz, M. R. (2014). False discovery rate control is a recommended alternative to Bonferroni-type adjustments in health studies. *Journal of Clinical Epidemiology*, **67**(8), 850–857.
- [31] Gualdi-Russo, E. and Toselli, S. (2002). Influence of various factors on the measurement of multifrequency bioimpedance. *Homo*, **53**(1), 1–16.
- [32] Hajika, R., Gunasekaran, T. S., Haigh, C. D. S. Y., Pai, Y. S., Hayashi, E., Lien, J., Lottridge, D., and Billingham, M. (2024). RadarHand: A wrist-worn radar for on-skin touch-based proprioceptive gestures. *ACM Transactions on Computer-Human Interaction*, **31**(2), 1–36.
- [33] Hanke, T. (2004). HamNoSys – representing sign language data in language resources and language processing contexts. In *International Conference on Language Resources and Evaluation (LREC)*, volume 4, pages 1–6.
- [34] Harrigan, J. A. (1985). Self-touching as an indicator of underlying affect and language processes. *Social Science & Medicine*, **20**(11), 1161–1168.
- [35] Harrison, C., Tan, D., and Morris, D. (2010). Skinput: Appropriating the body as an input surface. In *Proceedings of the ACM Conference on Human Factors in Computing Systems (CHI)*, pages 453–462, Atlanta, Georgia, USA. ACM.
- [36] Harrison, C., Benko, H., and Wilson, A. D. (2011). OmniTouch: Wearable multitouch interaction everywhere. In *Proceedings of the ACM Symposium on User Interface Software and Technology (UIST)*, pages 441–450, Santa Barbara, California, USA. ACM.
- [37] Hearle, J. W. S. (1952). The electrical resistance of textile materials: A review of the literature. *Journal of the Textile Institute Proceedings*, **43**(4), P194–P223.
- [38] Hodges, P. W. and Bui, B. H. (1996). A comparison of computer-based methods for the determination of onset of muscle contraction using electromyography. *Electroencephalography and Clinical Neurophysiology/Electromyography and Motor Control*, **101**(6), 511–519.
- [39] Hosseini, M., Ihmels, T., Chen, Z., Koelle, M., Müller, H., and Boll, S. (2023). Towards a consensus gesture set: A survey of mid-air gestures in HCI for maximized agreement across domains. In *Proceedings of the ACM Conference on Human Factors in Computing Systems (CHI)*, pages 1–24.
- [40] Huang, C.-H., Yi, H., Höschle, M., Safroshkin, M., Alexiadis, T., Polikovskiy, S., Scharstein, D., and Black, M. J. (2022). Capturing and inferring dense full-body human–scene contact. In *Computer Vision and Pattern Recognition (CVPR)*, pages 13274–13285.

- [41] Joo, H., Liu, H., Tan, L., Gui, L., Nabbe, B., Matthews, I., Kanade, T., Nobuhara, S., and Sheikh, Y. (2015). Panoptic studio: A massively multiview system for social motion capture. In *International Conference on Computer Vision (ICCV)*, pages 3334–3342.
- [42] Kanakri, S. M., Shepley, M., Varni, J. W., and Tassinary, L. G. (2017). Noise and autism spectrum disorder in children: An exploratory survey. *Research in Developmental Disabilities*, **63**, 85–94.
- [43] Khalil, S. F., Mohktar, M. S., and Ibrahim, F. (2014). The theory and fundamentals of bioimpedance analysis in clinical status monitoring and diagnosis of diseases. *Sensors*, **14**(6), 10895–10928.
- [44] Klima, E. S. and Bellugi, U. (1979). *The signs of language*. Harvard University Press.
- [45] Kocabas, M., Athanasiou, N., and Black, M. J. (2020). VIBE: Video inference for human body pose and shape estimation. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5252–5262, Piscataway, NJ.
- [46] Kratimenos, A., Pavlakos, G., and Maragos, P. (2021). Independent sign language recognition with 3D body, hands, and face reconstruction. In *International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 4270–4274.
- [47] Krishna, S., Vijay, V. P., and Dinesh, B. J. (2021). SignPose: Sign language animation through 3D pose lifting. In *International Conference on Computer Vision (ICCV)*, pages 2640–2649.
- [48] Kronrod, A. and Ackerman, J. M. (2019). I’m so touched! Self-touch increases attitude extremity via self-focused attention. *Acta Psychologica*, **195**, 12–21.
- [49] Kubo, Y., Koguchi, Y., Shizuki, B., Takahashi, S., and Hilliges, O. (2019). Audiotouch: Minimally invasive sensing of micro-gestures via active bio-acoustic sensing. In *Proceedings of the International Conference on Human-Computer Interaction with Mobile Devices and Services*, pages 1–13.
- [50] Kwok, Y. L. A., Gralton, J., and McLaws, M.-L. (2015). Face touching: A frequent habit that has implications for hand hygiene. *American Journal of Infection Control*, **43**(2), 112–114.
- [51] Kyle, U. G., Bosaeus, I., De Lorenzo, A. D., Deurenberg, P., Elia, M., Gómez, J. M., Heitmann, B. L., Kent-Smith, L., Melchior, J.-C., Pirlich, M., *et al.* (2004). Bioelectrical impedance analysis—part I: Review of principles and methods. *Clinical Nutrition*, **23**(5), 1226–1243.

- [52] Łacheta, J., Czajkowska-Kisil, M., Linde-Usiekiewicz, J., and Rutkowski, P. (2016). Korpusowy słownik polskiego języka migowego/Corpus-based Dictionary of Polish Sign Language.
- [53] Laput, G., Xiao, R., Chen, X. A., Hudson, S. E., and Harrison, C. (2014). Skin buttons: Cheap, small, low-powered and clickable fixed-icon laser projectors. In *Proceedings of the ACM Symposium on User Interface Software and Technology (UIST)*, pages 389–394, Honolulu, Hawaii, USA. ACM.
- [54] Li, T., Bolkart, T., Black, M. J., Li, H., and Romero, J. (2017). Learning a model of facial shape and expression from 4D scans. *ACM Transactions on Graphics (Proceedings of SIGGRAPH Asia)*, **36**(6), 194–1.
- [55] Loper, M., Mahmood, N., Romero, J., Pons-Moll, G., and Black, M. J. (2015). SMPL: A skinned multi-person linear model. *ACM Transactions on Graphics*, **34**(Article 248).
- [56] Mahmood, N., Ghorbani, N., F. Troje, N., Pons-Moll, G., and Black, M. J. (2019a). AMASS: Archive of motion capture as surface shapes. In *International Conference on Computer Vision (ICCV)*, pages 5441–5450.
- [57] Mahmood, N., Ghorbani, N., Troje, N. F., Pons-Moll, G., and Black, M. J. (2019b). AMASS: Archive of motion capture as surface shapes. In *Proceedings of the IEEE/CVF International Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5442–5451.
- [58] Majchrowska, S., Plantykov, M., and Olech, M. (2022). Handling sign language transcription system with the computer-friendly numerical multilabels.
- [59] Mathews, R. J. and Jovanov, E. (2023). Enabling complex impedance spectroscopy for cardio-respiratory monitoring with wearable biosensors: A case study. *Electrochem*, **4**(3), 389–410.
- [60] Matthie, J. R. (2008). Bioimpedance measurements of human body composition: critical analysis and outlook. *Expert Review of Medical Devices*, **5**(2), 239–261.
- [61] Meshcapade (2025). Meshcapade GmbH, Tübingen, Germany. <https://meshcapade.com>.
- [62] Mollyn, V. and Harrison, C. (2024). EgoTouch: On-body touch input using AR/VR headset cameras. In *Proceedings of the 37th Annual ACM Symposium on User Interface Software and Technology*, pages 1–11.
- [63] Moryossef, A., Tsochantaridis, I., Dinn, J., Camgoz, N. C., Bowden, R., Jiang, T., Rios, A., Muller, M., and Ebling, S. (2021). Evaluating the immediate applicability

- of pose estimation for sign language recognition. In *Computer Vision and Pattern Recognition (CVPR)*, pages 3434–3440.
- [64] Mueller, S. M., Martin, S., and Grunwald, M. (2019). Self-touch: Contact durations and point of touch of spontaneous facial self-touches differ depending on cognitive and emotional load. *PLOS ONE*, **14**(3), e0213677.
- [65] Mujibiya, A., Cao, X., Tan, D. S., Morris, D., Patel, S. N., and Rekimoto, J. (2013). The sound of touch: On-body touch and gesture sensing based on transdermal ultrasound propagation. In *Proceedings of the ACM International Conference on Interactive Tabletops and Surfaces*, pages 189–198, St. Andrews Scotland, United Kingdom. ACM.
- [66] Müller, L. (2024). *Self- and Interpersonal Contact in 3D Human Mesh Reconstruction*. PhD thesis, University of Tübingen, Tübingen, Germany.
- [67] Müller, L., Osman, A. A. A., Tang, S., Huang, C.-H. P., and Black, M. J. (2021). On self-contact and human pose. In *Computer Vision and Pattern Recognition (CVPR)*, pages 9990–9999.
- [68] Naranjo-Hernández, D., Reina-Tosina, J., and Min, M. (2019). Fundamentals, recent advances, and future challenges in bioimpedance devices for healthcare applications. *Journal of Sensors*, **2019**.
- [69] Nocedal, J. and Wright, S. J. (2006). *Numerical Optimization*. Springer, New York, NY, USA, second edition.
- [70] Pang, H. T., Canarlan, F., and Chu, M. (2022). Individual differences in conversational self-touch frequency correlate with state anxiety. *Journal of Nonverbal Behavior*, **46**(3), 299–319.
- [71] Pavlakos, G., Choutas, V., Ghorbani, N., Bolkart, T., Osman, A. A. A., Tzionas, D., and Black, M. J. (2019). Expressive body capture: 3D hands, face, and body from a single image. In *Computer Vision and Pattern Recognition (CVPR)*, pages 10975–10985.
- [72] Pfau, R., Steinbach, M., and Woll, B. (2012). *Sign language*. De Gruyter Mouton.
- [73] Porfirio, A. J., Wiggers, K. L., Oliveira, L. E., and Weingaertner, D. (2013). LIBRAS sign language hand configuration recognition based on 3D meshes. In *2013 IEEE International Conference on Systems, Man, and Cybernetics*, pages 1588–1593. IEEE.
- [74] Rahman, J., Mumin, J., and Fakhrudin, B. (2020). How frequently do we touch facial T-Zone: A systematic review. *Annals of Global Health*, **86**(1).

- [75] realsasl (2025). Real SASL: Real South African Sign Language. <https://www.realsasl.com/>.
- [76] Renz, K., Stache, N. C., Fox, N., Varol, G., and Albanie, S. (2021). Sign segmentation with changepoint-modulated pseudo-labelling. In *Computer Vision and Pattern Recognition Workshops (CVPRw)*.
- [77] Romero, J., Tzionas, D., and Black, M. J. (2017). Embodied hands: Modeling and capturing hands and bodies together. *ACM Transactions on Graphics (Proceedings of SIGGRAPH Asia)*, **36**(6), 1–17.
- [78] Rong, Y., Shiratori, T., and Joo, H. (2021). FrankMocap: A monocular 3D whole-body pose estimation system via regression and integration. In *International Conference on Computer Vision Workshops (ICCVw)*, pages 1749–1759.
- [79] Sandler, W. and Lillo-Martin, D. C. (2006). *Sign language and linguistic universals*. Cambridge University Press.
- [80] Sato, M., Poupyrev, I., and Harrison, C. (2012). Touché: Enhancing touch interaction on humans, screens, liquids, and everyday objects. In *Proceedings of the ACM Conference on Human Factors in Computing Systems (CHI)*, pages 483–492, Austin, Texas, USA. ACM.
- [81] Schielzeth, H., Dingemanse, N. J., Nakagawa, S., Westneat, D. F., Allogue, H., Teplitsky, C., Réale, D., Dochtermann, N. A., Garamszegi, L. Z., and Araya-Ajoy, Y. G. (2020). Robustness of linear mixed-effects models to violations of distributional assumptions. *Methods in Ecology and Evolution*, **11**(9), 1141–1152.
- [82] Shoemake, K. (1985). Animating rotation with quaternion curves. In *International Conference on Computer Graphics and Interactive Techniques (SIGGRAPH)*, pages 245–254.
- [83] Stokoe, W. C. (1978). *Sign Language Structure*. Linstok Press, Silver Spring, MD. Reprinted from original 1960 publication: *Studies in Linguistics, Occasional Papers*, Vol. 8, University of Buffalo.
- [84] Sun, Q., Wang, Y., Zeng, A., Yin, W., Wei, C., Wang, W., Mei, H., Leung, C.-S., Liu, Z., Yang, L., and Cai, Z. (2024). AiOS: All-in-one-stage expressive human pose and shape estimation. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1834–1843.
- [85] Tennant, R. A., Gluszak, M., and Brown, M. G. (2010). *The American Sign Language Handshape Dictionary*. Gallaudet University Press.

- [86] Theodorakis, S., Pitsikalis, V., and Maragos, P. (2014). Dynamic–static unsupervised sequentiality, statistical subunits and lexicon for sign language recognition. *Image and Vision Computing (IVS)*, **32**(8), 533–549.
- [87] U.S. Department of Justice (1990). Americans with Disabilities Act of 1990, as amended. <https://www.ada.gov/>. 42 U.S.C. § 12101 et seq.
- [88] van der Kooij, E. (1997). Contact: A phonological or a phonetic feature of signs? *Linguistics in the Netherlands*, **14**(1), 109–122.
- [89] Vázquez-Enríquez, M., Alba-Castro, J. L., Docío-Fernández, L., and Rodríguez-Banga, E. (2021). Isolated sign language recognition with multi-scale spatial-temporal graph convolutional networks. In *Computer Vision and Pattern Recognition (CVPR)*, pages 3462–3471.
- [90] Ward, L. C. (2019). Bioelectrical impedance analysis for body composition assessment: Reflections on accuracy, clinical utility, and standardisation. *European Journal of Clinical Nutrition*, **73**(2), 194–199.
- [91] World Health Organization (2010). A healthy lifestyle – WHO recommendations. <https://www.who.int/europe/news-room/fact-sheets/item/a-healthy-lifestyle---who-recommendations>.
- [92] World Health Organization (2025). Deafness and hearing loss. <https://www.who.int/news-room/fact-sheets/detail/deafness-and-hearing-loss>.
- [93] Xu, C., Solomon, S. A., and Gao, W. (2023). Artificial intelligence-powered electronic skin. *Nature Machine Intelligence*, **5**(12), 1344–1355.
- [94] Zhang, H., Tian, Y., Zhang, Y., Li, M., An, L., Sun, Z., and Liu, Y. (2023). PyMAF-X: Towards well-aligned full-body model regression from monocular images. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, **45**(10), 12287–12303.
- [95] Zhang, N., Jia, W., Wang, P., King, M.-F., Chan, P.-T., and Li, Y. (2020). Most self-touches are with the nondominant hand. *Scientific Reports*, **10**(1), 1–13.
- [96] Zhang, Y., Zhou, J., Laput, G., and Harrison, C. (2016). Skintrack: Using the body as an electrical waveguide for continuous finger tracking on the skin. In *Proceedings of the ACM Conference on Human Factors in Computing Systems (CHI)*, pages 1491–1503, San Jose, California, USA. ACM.
- [97] Zhang, Y., Kienzle, W., Ma, Y., Ng, S. S., Benko, H., and Harrison, C. (2019). Actitouch: Robust touch detection for on-skin AR/VR interfaces. In *Proceedings of the 32nd Annual ACM Symposium on User Interface Software and Technology*, pages 1151–1159.

## *Bibliography*

---

- [98] Zuo, R., Wei, F., Chen, Z., Mak, B., Yang, J., and Tong, X. (2024). A simple baseline for spoken language to sign language translation with 3D avatars. In *European Conference on Computer Vision*, pages 36–54. Springer.
- [99] Zurich Instruments AG (2022). *ziMFIA User Manual*. Accessed: 2024-09-30.