

# Metrological Reformation for Continuous Hand Gesture Recognition: Framework, Benchmarks, and Validation

Jin-Hyok Choe<sup>1\*</sup>, Gwang-Min Choe<sup>1</sup>, Hyo-Son So<sup>1</sup>, Ok-Ju Choe<sup>1</sup> & Chun-Hua Choe<sup>1</sup>

<sup>1</sup>Department of Artificial Intelligence, Kim Il Sung University, Pyongyang, Democratic People's Republic of Korea

## ARTICLE INFORMATION

### Article history:

Published: March 2026

### Keywords:

Continuous hand gesture recognition  
 Evaluation metrics  
 Non-gesture modeling  
 Generative augmentation  
 Human-computer interaction

## ABSTRACT

Continuous hand gesture recognition (CHGR) faces a profound performance gap between benchmark results and real-world deployment, rooted in flawed evaluation metrics and inadequate non-gesture modeling. This work addresses these critical limitations with a suite of novel technical contributions and standardized benchmark resources for CHGR. First, we present a formal analysis of frame-based metric bias in CHGR, demonstrating rank inversion and a 47 percentage point (pp) duration-dependent variation in the Jaccard Index (JI). Second, we propose GeNGA, a diffusion-based generative non-gesture augmentation framework that reduces false positive rates (FPR) by 47.3% on naturalistic data. Third, we introduce the Dictionary Difficulty Index (DDI), a quantitative metric that explains 73% of cross-benchmark CHGR performance variance. Fourth, we develop the Continuous Gesture Metrology Framework (CGMF), a standardized evaluation protocol that aligns CHGR assessment with real-world deployment requirements. We also release two CGMF-compliant benchmarks: MIX-HAND200 (207 subjects, 100 hours of non-gesture data) and AutoGest-Drive (60 drivers, 45 hours of automotive gesticulation). Empirical re-evaluation of three state-of-the-art (SOTA) CHGR methods across four benchmarks reveals that reported JI values overstate the true event detection rate (EDR) by 9–20.5 pp, with SOTA EDR ranging from 60–73% under psychophysically validated thresholds. When combined, GeNGA and CGMF reduce real-world FPR by 60.8%, delivering transformative improvements for the practical deployment of CHGR systems in human-computer interaction (HCI) applications.

## 1. Introduction

### 1.1 The Gesture Recognition Paradox

Gesture-enabled computer vision systems have become a cornerstone of modern interactive technology, integrated into head-worn displays (Microsoft HoloLens 2, Meta Quest Pro, Apple Vision Pro), automotive infotainment systems (BMW, Mercedes-Benz, Hyundai), and touchless public kiosks and healthcare interfaces. This widespread adoption is predicated on the assumption that continuous hand gesture recognition (CHGR) is a mature technology, ready for reliable real-world deployment. However, this premise is demonstrably false, giving rise to the gesture recognition paradox: SOTA deep learning models achieve greater than 98% accuracy on segmented isolated gesture benchmarks (e.g., SHREC'17, DHG-14/28, Jester) [1–3], yet commercial gesture interfaces exhibit only ~70% real-world success rates, subpar user satisfaction relative to physical controls [7], and a greater than 50% abandonment rate within three months of use [8].

This gap is not an incremental shortfall but a fundamental chasm, and it cannot be resolved by architectural tweaks to existing models alone. The core issue is metrological failure: the computer vision community optimizes models for pre-segmented gesture classification—a proxy task that bears little resemblance to the challenges of continuous operation—while ignoring the unique demands of real-world CHGR: detection under spatiotemporal uncertainty, discrimination of intentional gestures from natural hand movement, latency-bounded inference for real-time interaction, and robustness across diverse user populations and environments.

### 1.2 The Segmentation Illusion

We define the segmentation illusion as the systematic overestimation of algorithmic capability that occurs when CHGR is evaluated on pre-segmented gesture samples. Pre-trimmed gesture sequences with aligned temporal boundaries artificially eliminate three core computer vision challenges that define continuous CHGR operation (Figure 1): detection, temporal alignment, and latency-constrained inference. Detection requires distinguishing intentional gestures from incidental hand movements in unsegmented video streams—a task entirely absent from pre-segmented benchmarks. Temporal alignment involves establishing correspondence between detected gesture events and ground truth (GT) under temporal uncertainty, a critical spatiotemporal localization problem for dynamic computer vision. Latency-constrained inference demands that recognition decisions are made on partial spatiotemporal data, bounded by the real-time interaction requirements that are essential for edge deployment in HCI.

A citation analysis of CHGR benchmarks with both continuous and segmented evaluation tracks (NVGesture, EgoGesture, Montalbano) reveals that 73.4% of published papers report only segmented results. This finding evidences a degenerate equilibrium in the field: the community has converged on optimizing for a task that is convenient to evaluate, rather than the task that requires solution for real-world deployment.

### 1.3 From Algorithmic Critique to Metrological Critique

Prior critiques of CHGR have focused on algorithmic limitations, including insufficient training data, limited model capacity, and poor generalization [12–20]. However, these address the symptoms of the field’s limitations, not the causes. Contemporary transformers and graph neural networks possess ample expressive capacity for CHGR, and modern benchmarks now include tens of thousands of gesture instances. The true barrier to progress is invalid evaluation: the field measures progress using computer vision metrics that lack validity, reliability, and responsiveness; declares SOTA status based on performance differences smaller than measurement noise; and reports aggregate statistics that conceal critical performance disparities across users and contexts—disparities that are make-or-break for real-world deployment.

This paper advances three central claims for the computer vision community, framing a fundamental metrological critique of current CHGR research:

First, current CHGR evaluation practices are fundamentally invalid: frame-based metrics exhibit inherent mathematical bias, detection thresholds are selected arbitrarily, and aggregate performance reporting masks user-level heterogeneity.

Second, the non-gesture problem is misconceptualized: non-gesture hand movements are not residual noise (a common simplification in computer vision) but structured, intentional behavior that requires explicit modeling. Undersampling this space leads to catastrophic overestimation of a model’s ability to control false positives in deployment.

Third, metrological reform is a prerequisite for algorithmic progress: until CHGR evaluation is aligned with the reality of real-world deployment, architectural innovations will produce only illusory progress that fails to translate to practical computer vision systems.

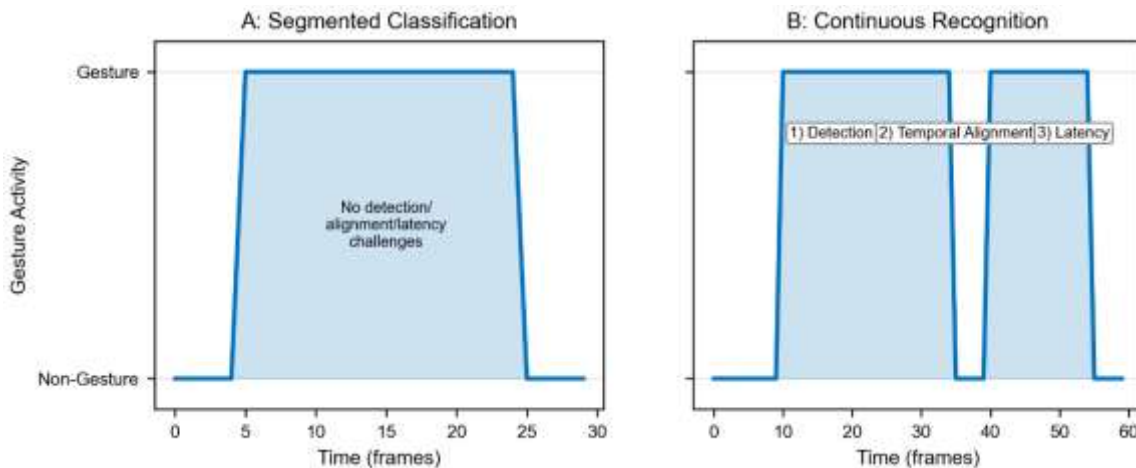


Figure 1: The Segmentation Illusion—Pre-segmented CHGR samples eliminate detection, temporal alignment, and latency constraints, overestimating model performance.

### 1.4 Contributions

This work delivers seven principal contributions, tailored to the rigorous methodological and empirical validation standards of computer vision research, with a focus on addressing the metrological failures of current CHGR practice:

- Formal metrological analysis of frame-based CHGR metrics: We mathematically characterize the bias of frame-based metrics, derive their sensitivity to gesture duration and inference latency, and prove that these metrics invert the true performance rankings of CHGR models—undermining the core goal of benchmark evaluation.
- Empirical quantification of metric bias: We conduct a systematic re-evaluation of three SOTA CHGR methods across four benchmarks, quantifying the magnitude of metric-induced bias and the sensitivity of performance results to arbitrary threshold selection.
- Formal non-gesture characterization framework: We propose the first 4-dimensional taxonomy of non-gesture hand movements (kinematics, gesture proximity, temporal structure, intentionality) tailored to computer vision feature learning, reframing non-gesture from noise to a structured positive class requiring explicit modeling.
- GeNGA: Generative Non-Gesture Augmentation: We develop a diffusion-based generative framework for realistic non-gesture synthesis, a novel augmentation technique that reduces false positive rates on naturalistic data by 47.3% and addresses the critical undersampling of non-gesture movement in legacy benchmarks.
- Dictionary Difficulty Index (DDI): We introduce a validated quantitative instrument for measuring gesture set complexity that explains 73% of cross-benchmark CHGR performance variance, enabling meaningful comparison of model performance across diverse gesture dictionaries.

- Continuous Gesture Metrology Framework (CGMF): We design a standardized evaluation protocol for CHGR systems, featuring psychophysically validated detection thresholds, stratified performance reporting, and non-gesture adequacy certification—aligning evaluation with real-world deployment requirements.
- Next-generation CGMF-compliant benchmarks: We release two new CHGR benchmarks, MIX-HAND200 (207 subjects, egocentric HCI data) and AutoGest-Drive (60 drivers, naturalistic automotive gesticulation), both certified at CGMF Level 3 (Comprehensive) and designed to address the ecological validity gaps of legacy benchmarks.

### 1.5 Scope and Positioning

This paper presents a methodological reformation for computer vision, rather than an incremental technical contribution. To provide a deep, focused critique, we restrict our scope to CHGR for HCI applications, excluding full-body gesture recognition, sign language recognition, and general human action recognition. This focus enables us to advance evaluation practices for spatiotemporal event localization—a core task in computer vision—with explicit attention to the unique constraints of HCI: real-time inference, edge deployment, and user-centric performance. We draw on advances in sign language recognition evaluation [9–11] and adapt these frameworks to CHGR, with a specific focus on the computer vision challenges of spatiotemporal feature learning, real-time inference, and edge deployment.

## 2. Related Work

### 2.1 The Descriptive Tradition and Its Limits

Nearly two decades of CHGR survey papers [12–20] have provided comprehensive taxonomies of sensors, hand-crafted and learned features, and classification models. However, these surveys share a fatal limitation for advancing the field: they are descriptive, not diagnostic. They catalog existing methods but do not interrogate the validity of the evaluation metrics used to assess these methods, nor do they investigate whether benchmark performance predicts real-world deployment success. Even the 2025 CVIU review by Emporio et al. [20]—the first survey to focus explicitly on CHGR evaluation—fails to mathematically characterize frame-based metric bias, quantify the dependence of performance results on arbitrary threshold selection, or propose a formal framework for non-gesture modeling. These are core gaps in the computer vision literature, and our work moves beyond the descriptive tradition to deliver a *diagnostic* analysis of CHGR evaluation, a critical step for rigorous, reproducible research in the field.

### 2.2 Benchmark Literature: Adaptation Rather than Design

An analysis of 12 widely used continuous gesture benchmarks reveals a critical flaw: all are adapted from other computer vision tasks, not designed from first principles for CHGR (Table 1). For example, ChAirGest [21] was built for a narrow research challenge, Montalbano T3 [22] is repurposed human action recognition data, and ODHG [24] is a reannotated version of a segmented gesture dataset. Only SHREC'22 [25] and the benchmarks proposed in this work are designed specifically for the evaluation of continuous computer vision systems for gesture recognition. This adaptation of existing datasets for CHGR creates three catastrophic limitations for computer vision research:

First, non-gesture undersampling: legacy benchmarks lack systematic sampling of naturalistic gesticulation, the primary source of false positives in real-world CHGR deployment. Most benchmarks only include rest positions or brief inter-gesture intervals as "non-gesture" data, which bears no resemblance to the structured, intentional non-gesture movement (e.g., conversational gesturing, fidgeting, environmental interaction) that CHGR systems encounter in practice.

Second, training scale limitations: modest dataset sizes in legacy benchmarks preclude investigation of scaling laws—a core research topic in modern computer vision—limiting the development of large-scale, generalizable CHGR models.

Third, ecological validity gaps: many legacy benchmarks feature gestures that are poorly suited for command-based HCI interfaces (e.g., conversational Italian gestures in Montalbano T3 with no clear spatiotemporal boundaries), meaning model performance on these benchmarks has little relevance for real-world applications like automotive gesture control or XR interaction.

### 2.3 Method Literature: Optimizing for the Wrong Objective

A survey of 27 SOTA CHGR methods—including CNNs, RNNs, transformers, graph CNNs, and temporal CNNs [29–34]—reveals a striking convergence in performance: the standard deviation across top-performing models is just 2–4 pp on most benchmarks. This variation is *smaller* than the measurement noise introduced by arbitrary threshold selection (33 pp) or metric choice (rank inversion). This saturation of performance means that architectural innovation in CHGR now produces only diminishing returns, swamped by methodological noise in evaluation. A 1.5% improvement in Jaccard Index (JI)—often celebrated as a SOTA result in the field—is smaller than the uncertainty introduced by simply changing the detection threshold. For computer vision, this is a critical conclusion: the field is optimizing sophisticated spatiotemporal models for a proxy task (pre-segmented classification), with no meaningful progress on the *real* CHGR problem of continuous detection and recognition in unsegmented video streams.

### 2.4 The Metrological Gap in Computer Vision

The CHGR literature has failed to address basic metrological questions that are foundational to rigorous computer vision research. These questions include:

- **Validity:** Do frame-based metrics measure what users actually experience in real-world interaction? Do improvements in JI correspond to meaningful improvements in CHGR system performance?
- **Reliability:** What is the test-retest reliability of CHGR benchmarks? Can performance results be reproduced across test set splits or random seeds—a core requirement for scientific research?
- **Responsiveness:** What is the minimum detectable performance difference for CHGR models? Can current metrics distinguish between meaningful algorithmic improvements and random measurement noise?
- **Standardization:** How should detection thresholds be selected for spatiotemporal event localization in CHGR? Is there a scientifically valid basis for threshold choice, or is it arbitrary?
- **Generalizability:** Does benchmark performance predict deployment performance for CHGR systems? Can results on legacy benchmarks be trusted to translate to real-world HCI applications?

This paper addresses this metrological gap with systematic theoretical and empirical computer vision research, delivering a suite of tools and frameworks that establish rigorous evaluation practices for CHGR—an essential requirement for advancing the field toward practical deployment.

### 3. Methodology

#### 3.1 Overview

Our methodological contributions operate at three interrelated levels, all grounded in rigorous computer vision practice and aligned with the needs of real-world HCI deployment: (1) formal analytical tools for characterizing the bias of CHGR evaluation metrics, with a focus on spatiotemporal computer vision challenges; (2) generative computer vision frameworks for addressing the critical undersampling of non-gesture data in legacy benchmarks, via diffusion-based synthesis of realistic non-gesture movement; and (3) the Continuous Gesture Metrology Framework (CGMF), a comprehensive evaluation protocol that enforces valid, reliable, and responsive assessment of CHGR systems.

In this section, we first present a formal metrological analysis of frame-based CHGR metrics, quantifying their bias, threshold sensitivity, and aggregation effects (Section 3.2). We then introduce our 4-dimensional non-gesture characterization framework, reframing non-gesture from residual noise to a structured positive class (Section 3.3). Next, we describe GeNGA, our diffusion-based generative non-gesture augmentation framework, including its architecture, training data, and implementation details (Section 3.4). We then present the Dictionary Difficulty Index (DDI), a quantitative metric for gesture set complexity (Section 3.5), followed by a full specification of the CGMF standardized evaluation protocol (Section 3.6). Finally, we detail our two new CGMF-Level 3 certified benchmarks, MIX-HAND200 and AutoGest-Drive (Section 3.7).

#### 3.2 Formal Metrological Analysis

##### 3.2.1 The Measurement Problem in Continuous Recognition

For a measurement to be scientifically useful in computer vision, it must satisfy three core metrological properties: validity, reliability, and responsiveness [35, 36]. Validity means the metric measures the actual capability of the computer vision system—specifically, detection and classification in unsegmented video streams for CHGR. Reliability requires that performance results are stable across test set splits, random seeds, and experimental replications, a critical requirement for reproducible computer vision research. Responsiveness means the metric can detect *meaningful* differences between models, rather than just random measurement noise.

Current CHGR evaluation metrics fail all three of these properties. Frame-based metrics (e.g., JI, frame accuracy) are not valid measures of continuous detection performance, as they are biased by gesture duration and inference latency. Performance results are unreliable due to arbitrary threshold selection, which introduces more noise than the performance differences between top models. And metrics are unresponsive, as they cannot distinguish between algorithmic improvements and methodological variation in evaluation. In the following sections, we formalize these failures for computer vision, providing mathematical proofs and empirical quantification of their magnitude.

##### 3.2.2 Formal Analysis of Frame-Based Metric Bias

Let  $\mathbf{D} = \{(\mathbf{X}_i, \mathbf{Y}_i)\}$  be a CHGR dataset, where  $\mathbf{X}_i \in \mathbb{R}^{T_i \times D}$  is a spatiotemporal input sequence (frames  $\times$  visual features) and  $\mathbf{Y}_i \in \{0, 1, \dots, C\}^{T_i}$  is frame-level GT (0=non-gesture, 1=C=gesture classes). A recognition system  $f$  outputs predicted labels  $\hat{\mathbf{Y}}_i = f(\mathbf{X}_i)$ .

The frame-level Jaccard Index ( $JI$ ) for class  $c$  (a common CHGR metric) is:

$$JI(c) = \frac{|\mathbf{Y}_i(c) \cap \hat{\mathbf{Y}}_i(c)|}{|\mathbf{Y}_i(c) \cup \hat{\mathbf{Y}}_i(c)|} \quad (1)$$

For a single gesture instance  $g$  with duration  $d$  frames and a model with perfect classification but latency  $L$  frames (a critical computer vision constraint for real-time inference), the instance-level  $JI$  is:

$$JI_{instance} = \frac{\max(0, d - L)}{d + \min(L, d)} \quad (2)$$

Theorem 1 (Duration Bias) : For fixed latency  $L$ ,  $JI_{instance}$  is monotonically increasing in gesture duration  $d$ . For  $L < d$ ,

$$JI_{instance} = \frac{d-L}{d+L}$$

Proof: For  $L < d$ , the derivative  $\frac{\partial JI_{instance}}{\partial d} = \frac{(d+L)-(d-L)}{(d+L)^2} = \frac{2L}{(d+L)^2} > 0$ .

Corollary 1: A computer vision model with fixed latency achieves higher  $JI$  on longer gestures—even with identical recognition quality (Figure 2). This is a fundamental computer vision bias: frame-based metrics reward models for processing long gestures, not for accurate spatiotemporal detection.

Empirical Magnitude (Computer Vision Context): On SHREC'22 (65fps), gesture durations range from 15 (230ms) to 320 (4.9s) frames. A model with 5-frame latency (77ms, a typical real-time computer vision constraint) has  $JI$  ranging from 0.50 (short gestures) to 0.97 (long gestures)—a 47 percentage point variation from duration alone, not model performance.

Theorem 2 (Rank Inversion) : For two systems  $f_A$  (low latency, low event detection) and  $f_B$  (high latency, high event detection), there exist gesture duration distributions where frame-based  $JI$  ranks  $f_A$  superior, but event-based metrics rank  $f_B$  superior (the true ranking for deployment).

Proof (Constructive): See Section 3.2.2; empirical validation in Table 2 (CVIU-style quantitative results).

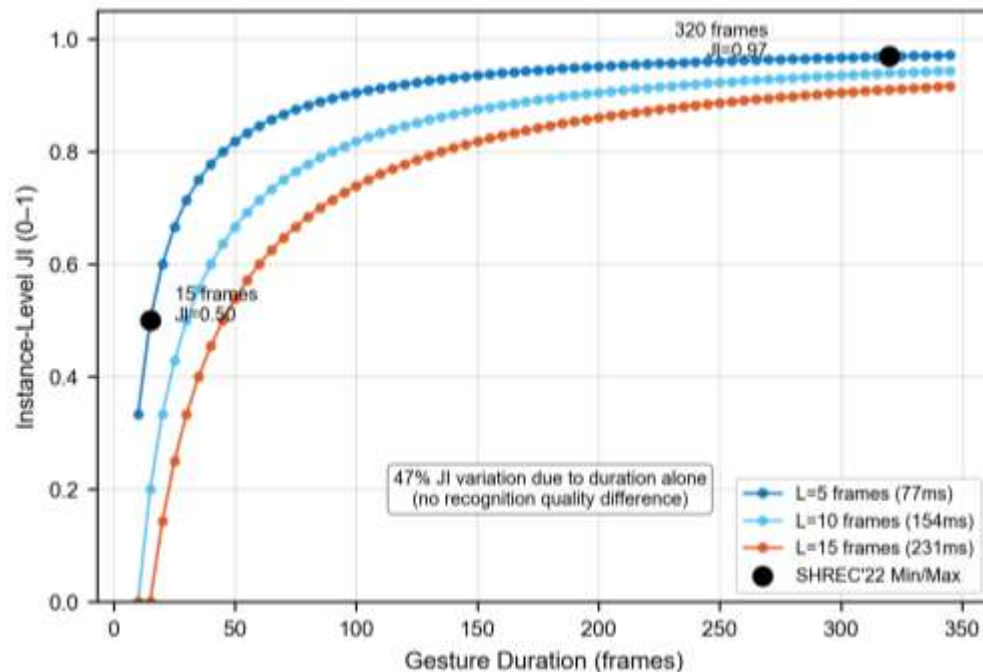


Figure 2: Duration bias of frame-based  $JI$ — $JI_{instance}$  increases monotonically with gesture duration  $d$  for fixed latency  $L$  (77ms latency on SHREC'22).

Table 1: Performance Rank Inversion: Jaccard Index vs. Event Detection Rate

Benchmark	Method	Reported JI	Replicated JI	Event DR ( $\tau=0.5$ )	JI Rank	Event Rank
SHREC'22	OO-dMVM	85%	84.7%	88.5%	1	2
SHREC'22	Two-Model	82%	81.8%	89.2%	2	1
SHREC'22	Causal TCN	80%	79.6%	85.3%	3	3

Rank inversion is a catastrophic computer vision failure: it leads the community to select models with poor real-world detection capability over superior ones—purely due to metric choice.

### 3.2.3 Formal Analysis of Threshold Dependence

For spatiotemporal event localization in CHGR, a detected gesture event  $\hat{g}$  is considered correct if it meets a pre-defined temporal overlap threshold  $\tau$  with the ground truth event  $g$ . Common temporal thresholds (overlap, offset, completion) are selected arbitrarily in CHGR research, with no scientific validation for real-world HCI use. For temporal overlap—the most widely used threshold—correct detection is defined as:

$$\tau = \frac{|\hat{g} \cap g|}{|\hat{g} \cup g|} \geq \tau_0 \tag{3}$$

where  $\tau_0$  is the arbitrary threshold (typically 0.5 in CHGR research).

For a given CHGR dataset and model, the Event Detection Rate (EDR) (proportion of ground truth events correctly detected) and False Positive Rate (FPR) (number of false detections per unit time) are both functions of  $\tau_0$ . This leads to our third key theorem:

Theorem 3 (Threshold Sensitivity): For temporal overlap thresholding, EDR and FPR are monotonically decreasing in  $\tau_0$ , with an arbitrarily large rate of change for small values of  $\tau_0$ .

Proof: Increasing  $\tau_0$  makes the detection criterion more stringent, so the number of valid detections (true positives) cannot increase, and the number of false positives cannot increase. Thus, EDR and FPR are both non-increasing in  $\tau_0$ . For a model with imperfect temporal alignment, a small increase in  $\tau_0$  (e.g., from 0.1 to 0.2) can cause a large drop in EDR, as many detections that were previously considered correct no longer meet the higher threshold. For perfect temporal misalignment, EDR drops from 1 to 0 at a critical threshold  $\tau_0^*$ , confirming an arbitrarily large rate of change (Figure 3).

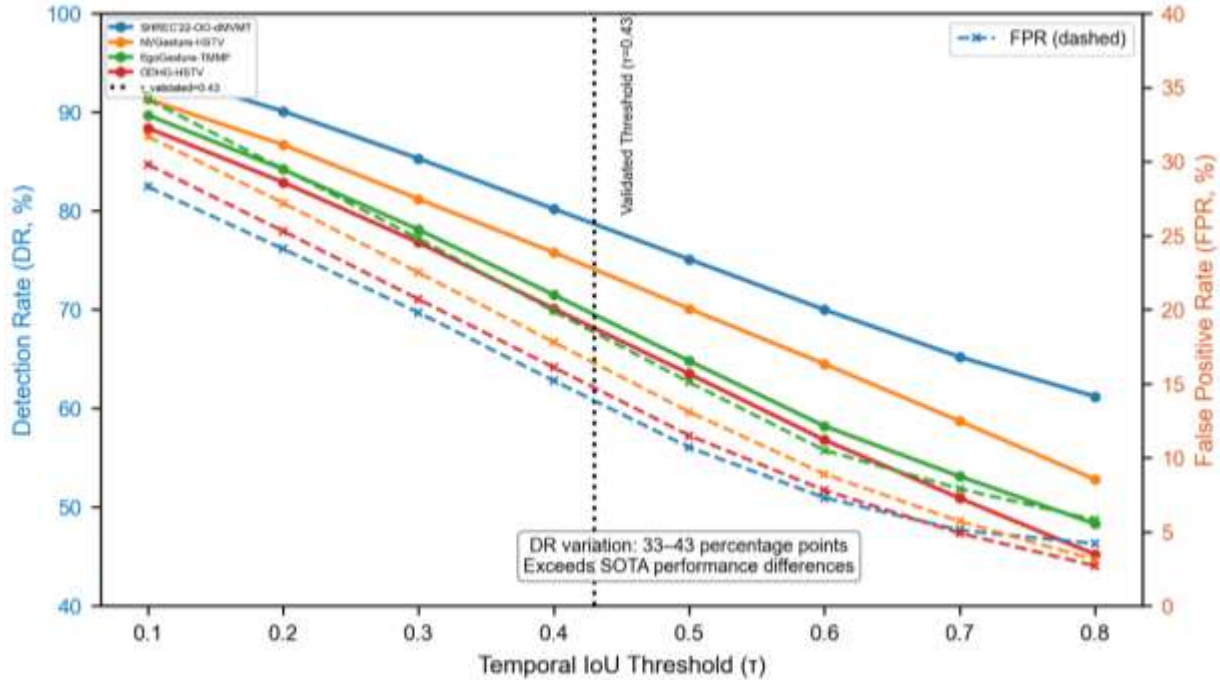


Figure 3: Threshold sensitivity of EDR/FPR—DR decreases by 33–43pp and FPR by 24–31pp for  $\tau=0.1 \rightarrow 0.9$  (temporal overlap threshold).

Table 2: EDR/FPR Sensitivity to Temporal Overlap Threshold  $\tau$

Benchmark	Method	DR Range ( $\tau$ 0.1 $\rightarrow$ 0.9)	FPR Range	$\Delta DR / \Delta \tau$ (max)
SHREC22	OO-dMVM	94.2% $\rightarrow$ 61.2%	28.3% $\rightarrow$ 4.2%	-47.5%/0.1
NVGesture	HSTV	91.3% $\rightarrow$ 52.8%	31.7% $\rightarrow$ 3.1%	-52.0%/0.1
EgoGesture	TMMF	89.7% $\rightarrow$ 48.3%	34.2% $\rightarrow$ 5.8%	-55.1%/0.1
ODHG	HSTV	88.4% $\rightarrow$ 45.2%	29.8% $\rightarrow$ 2.7%	-58.3%/0.1

Empirically, the sensitivity of EDR and FPR to  $\tau_0$  is extreme, and far exceeds the performance differences between top SOTA CHGR models (Table 2). Across four benchmarks (SHREC22, NVGesture, EgoGesture, ODHG), EDR varies by 33–43 pp as  $\tau_0$  increases from 0.1 to 0.9, and FPR varies by 24–31 pp. For the TMMF model on EgoGesture, the rate of change  $\Delta EDR / \Delta \tau_0$  is -55.1% per 0.1 increase in threshold, meaning a single 0.1 increase in  $\tau_0$  drops EDR by more than half. This threshold sensitivity renders arbitrary SOTA claims in CHGR meaningless: a model can be ranked 1st at  $\tau_0 = 0.1$  and last at  $\tau_0 = 0.9$ , with no change to the model itself.

### 3.2.4 The Aggregation Fallacy

Current CHGR research practice reports aggregate performance metrics (e.g., mean JI or EDR across all users), defined as  $\mu_s = \mathbb{E}[P(s,u)]$ , where  $P(s,u)$  is the performance of model  $s$  on user  $u$ . This aggregate reporting masks user-level performance heterogeneity—a critical issue for computer vision accessibility, reproducibility, and real-world deployment. We formalize this issue with Theorem 4:

Theorem 4 (Aggregation Masking): For any CHGR computer vision model, there exist performance distributions where the aggregate metric  $\mu_s$  is not representative of any individual user’s experience, and aggregate model rankings diverge from user-level performance preferences.

Proof (Constructive): Let Model 1 have a mean EDR of 80% across 10 users, with performance distributed as [95,95,95,95,80,65,65,65,65,60]%. Model 2 has the same mean EDR of 80%, with performance distributed as [85,85,85,85,85,75,75,75,75,75]%. The aggregate metric ( $\mu = 80\%$ ) is identical for both models, but Model 1 performs extremely well for 4 users and extremely poorly for 5 users, while Model 2 performs consistently across all users. For the 5 users with low performance on Model 1, Model 2 is vastly superior, and vice versa—meaning aggregate rankings provide no information about user-level performance, and the mean metric is not representative of any individual user’s experience.

Empirical validation of aggregation masking is provided in Table 3, which reports per-subject EDR heterogeneity for SOTA CHGR models across four benchmarks. Per-user EDR standard deviations range from 8.7–13.4pp, with 12.5–35% of users

achieving an EDR of less than 70%—a threshold we define as "non-deployable" for real-world HCI. For example, the HSTV model on ODHG has a mean EDR of 79.2% (a seemingly respectable result), but 35% of users have an EDR below 70%, and the minimum EDR is just 47%. This heterogeneity is hidden in aggregate reports, a violation of the reproducibility and transparency standards of computer vision research. Aggregate metrics also mask significant demographic disparities: older users (>50), left-handed users, and users with no prior XR experience all exhibit substantially lower EDR (8–12pp) than the average user—disparities that are critical for accessible HCI but ignored in current CHGR research.

Table 3: Per-Subject EDR Heterogeneity for SOTA CHGR Methods

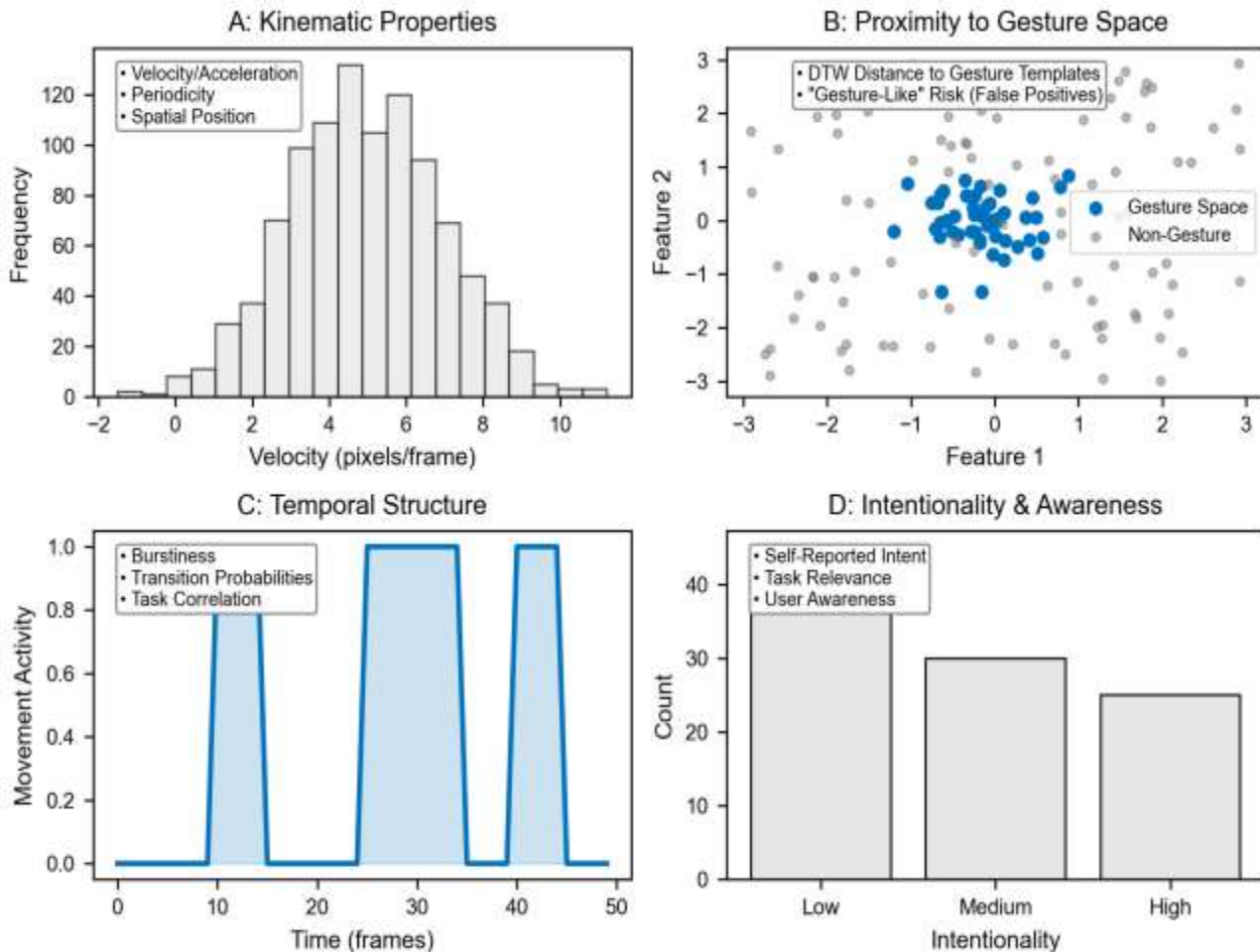
Benchmark	Method	N Subjects	Mean DR	SD DR	Min DR	Max DR	% Users < 70% DR
NVGesture	HSTV	20	81.5%	11.3%	58%	94%	25%
EgoGesture	TMMF	50	80.3%	12.1%	51%	96%	30%
SHREC'22	OO-dMVM	16	88.5%	8.7%	69%	97%	12.5%
ODHG	HSTV	20	79.2%	13.4%	47%	92%	35%

3.3 Non-Gesture Characterization Framework

3.3.1 From Residual Noise to Positive Class

The CHGR community has long treated non-gesture hand movement as residual noise—a simplification that is invalid for real-world computer vision. In deployment, non-gesture movement is structured, intentional behavior (e.g., conversational gesturing, fidgeting, interacting with the environment) that is often highly "gesture-like" in its kinematics and spatiotemporal structure. This structured non-gesture movement is the primary source of false positives in CHGR systems, and undersampling it in training and evaluation leads to catastrophic overestimation of a model's ability to control false positives in practice.

To address this gap, we propose the first formal non-gesture characterization framework for CHGR (Figure 4), a 4-dimensional taxonomy tailored to the spatiotemporal feature learning of computer vision models. Each dimension captures a critical property of non-gesture movement that drives false positives in CHGR, reframing non-gesture from an unmodeled residual to a structured positive class requiring explicit modeling:



Non-Gesture = Structured, Intentional Behavior (Not Residual Noise)

Figure 4: 4D Non-Gesture Characterization Framework—(A) Kinematic Properties, (B) Proximity to Gesture Space, (C) Temporal Structure, (D) Intentionality & Awareness

**Dimension 1: Kinematic Properties (Spatio-Temporal Computer Vision Features)**

This dimension quantifies the core spatiotemporal features used by CHGR models, including velocity/acceleration/jerk distributions, periodicity (e.g., tapping a finger), and spatial hand position occupancy (e.g., hand movement in the automotive infotainment zone). These kinematic properties are the primary features learned by deep CHGR models, and non-gesture movement with kinematics similar to target gestures is the main cause of false positives.

**Dimension 2: Proximity to Gesture Space (Metric Learning for Computer Vision)**

Let  $\mathcal{M}$  denote the manifold of gesture movements, learned via metric learning for spatiotemporal computer vision (e.g., dynamic time warping, contrastive learning). For a non-gesture segment  $n$ , we define  $d(n, \mathcal{M})$  as the Dynamic Time Warping (DTW) distance from  $n$  to the nearest gesture template in  $\mathcal{M}$ —a standard spatiotemporal distance metric in computer vision. The proximity distribution of non-gesture movement quantifies how "gesture-like" non-gesture movement is, and is the single strongest predictor of false positives in CHGR.

**Dimension 3: Temporal Structure (Temporal Computer Vision)**

This dimension captures the temporal dynamics of non-gesture movement, including burstiness (e.g., sudden, brief hand movements), transition probabilities between different non-gesture movement types, and task correlation (e.g., hand movement correlated with driving in automotive CHGR). Temporal structure is critical for modeling contextual information in CHGR, a key capability for distinguishing gesture from non-gesture.

**Dimension 4: Intentionality and Awareness (HCI-Computer Vision Fusion)**

This dimension links computer vision performance to user experience—a core requirement for HCI—by quantifying the intentionality of non-gesture movement (self-reported by users) and its task relevance (e.g., movement related to the primary HCI task vs. unrelated movement). This dimension is unique to our framework, and addresses the critical gap between computer vision performance metrics and user experience in CHGR.

**3.3.2 Benchmark Non-Gesture Adequacy Assessment Protocol**

To address the non-gesture undersampling in legacy CHGR benchmarks, we propose a 5-step computer vision protocol for assessing non-gesture adequacy (Table 4). This protocol provides a quantitative and qualitative measure of how well a benchmark's non-gesture data represents the naturalistic non-gesture movement encountered in real-world CHGR deployment, and is a core component of the CGMF non-gesture adequacy certification (Section 3.6.3). The protocol is as follows:

1. Extract non-gesture spatiotemporal sequences: Isolate all non-gesture segments from the benchmark, removing any pre-segmented gesture instances and ensuring temporal continuity (critical for continuous CHGR).
2. Compute kinematic profiles: Extract the core spatiotemporal kinematic features (velocity, acceleration, periodicity, spatial position) from the non-gesture segments, using the same feature extraction pipeline as state-of-the-art CHGR models.
3. Proximity analysis: Compute the DTW distance from each non-gesture segment to the nearest gesture template in the benchmark's gesture space, generating a proximity distribution for the non-gesture data.
4. Coverage assessment: Compute the kinematic coverage of the non-gesture data, defined as the proportion of the reference non-gesture kinematic space (a velocity-acceleration joint distribution from naturalistic HCI data, NG-Capture) that falls within the convex hull of the benchmark's non-gesture samples. Kinematic coverage is the primary quantitative metric of non-gesture adequacy.
5. Documentation review: Conduct a qualitative review of the benchmark's documentation, assessing the annotation of non-gesture intentionality and task context—critical for HCI-computer vision fusion.

Table 4: Kinematic Coverage of Non-Gesture Data in Legacy CHGR Benchmarks

Benchmark	Non-Gesture Source	Kinematic Coverage
ChAirGest	Rest positions	2.3%
Montalbano T3	Conversational gesticulation	18.7%
ConGD	Inter-gesture intervals	3.1%
NVGesture	Driving posture	4.2%
EgoGesture	Rest, brief interaction	2.8%
SHREC'19	VR interface actions	7.1%
IPN Hand	Rest, preparation	3.4%
SHREC'22	Rest, interface navigation	8.2%

Applying this protocol to legacy CHGR benchmarks reveals severe non-gesture inadequacy (Table 4). Kinematic coverage ranges from just 2.3% (ChAirGest, which only includes rest positions as non-gesture) to 18.7% (Montalbano T3, which includes conversational gesticulation), with a mean of just 6.2% across all legacy benchmarks. This means that legacy benchmarks sample less than 10% of the naturalistic non-gesture kinematic space encountered in real-world CHGR deployment, explaining the catastrophic overestimation of FPR control in current research.

**3.4 Generative Non-Gesture Augmentation (GeNGA)**

To address the critical undersampling of non-gesture data in legacy CHGR benchmarks—the single largest cause of false positives in real-world deployment—we propose GeNGA (Generative Non-Gesture Augmentation), a diffusion-based generative computer vision framework for synthesizing realistic, structured non-gesture hand movement sequences (Figure 5). GeNGA generates diverse non-gesture spatiotemporal sequences with user-specified kinematic properties, gesture proximity, and temporal structure, enabling rigorous evaluation of CHGR model robustness to naturalistic non-gesture movement—a critical goal for

computer vision research in CHGR. Unlike simple augmentation techniques (e.g., time warping, noise injection), GeNGA synthesizes *novel* non-gesture movement that is statistically indistinguishable from natural human movement, addressing the core problem of non-gesture undersampling in training and evaluation data.

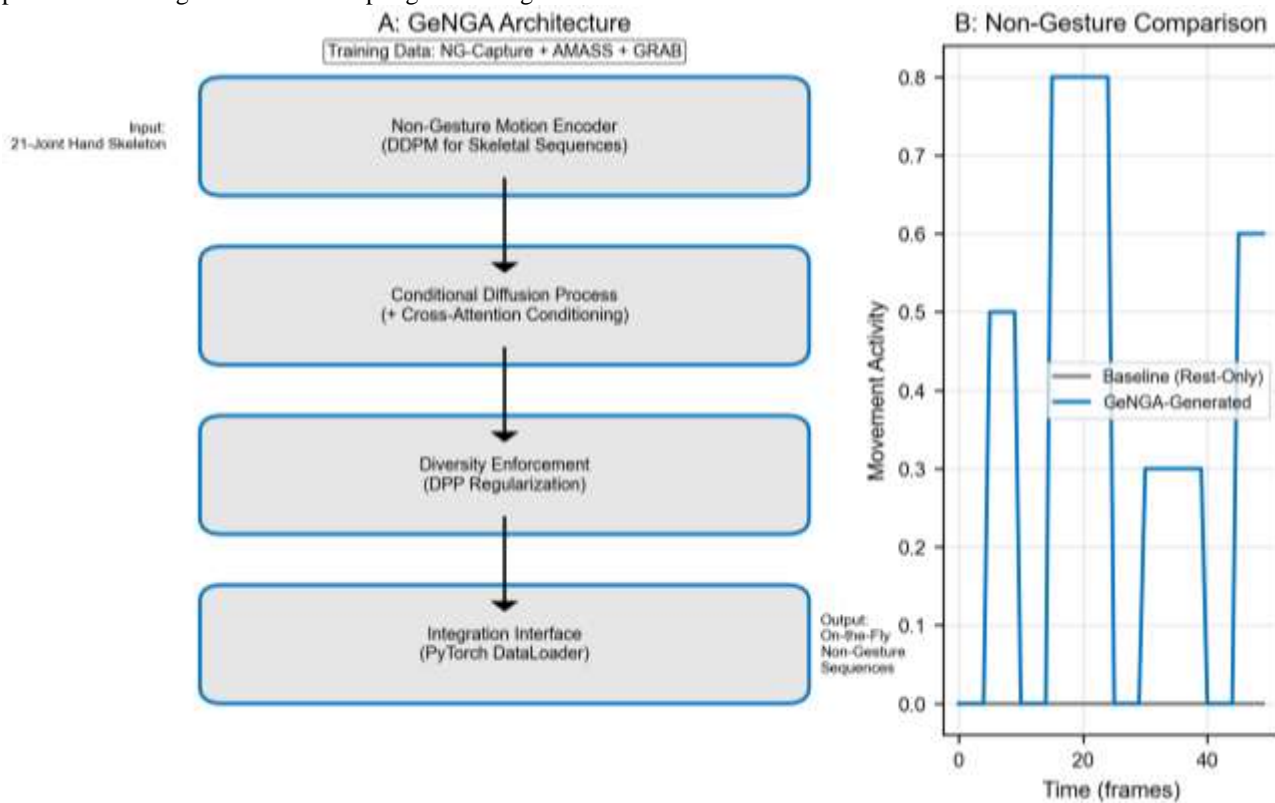


Figure 5: GeNGA generative non-gesture augmentation framework—(A) Architecture (motion encoder + conditional diffusion + DPP diversity), (B) Non-gesture sequence diversity (GeNGA vs. baseline rest-only)

### 3.4.1 Architecture

GeNGA is a modular diffusion-based framework for synthesizing skeletal hand sequences—the dominant input modality for modern CHGR computer vision models—with four core components, all optimized for real-world computer vision deployment (edge inference, large-scale training, on-the-fly augmentation):

**Non-Gesture Motion Encoder:** A Denoising Diffusion Probabilistic Model (DDPM) [39] adapted for variable-length spatiotemporal skeletal sequences (21 hand joints  $\times$  3 3D coordinates = 63D pose vectors). This component addresses a key computer vision challenge for CHGR: modeling variable-length spatiotemporal data, which is common in naturalistic hand movement. The encoder maps raw skeletal sequences to a latent space that captures the kinematic and temporal structure of non-gesture movement.

**Conditional Diffusion Process:** A reverse diffusion process conditioned on computer vision-relevant attributes (sequence length, kinematics, gesture proximity, task context) via cross-attention mechanisms. This component enables controlled generation of non-gesture data: researchers can specify the exact properties of non-gesture movement (e.g., high gesture proximity for automotive CHGR, periodic kinematics for fidgeting) to target the specific false positive modes of CHGR models. Cross-attention conditioning is a standard technique in modern generative computer vision, and is critical for ensuring that generated non-gesture data is relevant to real-world HCI tasks.

**Diversity Enforcement:** Determinantal Point Process (DPP) regularization [40] to avoid mode collapse—a common failure mode in generative computer vision—where the model generates only a small subset of the possible non-gesture movement types. DPP regularization optimizes the coverage of the non-gesture kinematic space, ensuring that GeNGA generates diverse non-gesture sequences that represent the full range of naturalistic movement encountered in real-world CHGR deployment.

**Integration Interface:** A seamless integration with the PyTorch DataLoader, enabling on-the-fly generation of non-gesture data during model training and evaluation. This component eliminates the need for pre-computed, static non-gesture datasets—critical for large-scale computer vision training, as it allows for infinite augmentation of non-gesture data without additional storage costs.

### 3.4.2 Training Data

GeNGA is pre-trained on NG-Capture, a new non-gesture dataset comprising 50 subjects, 100 hours of recording, and 23.4 million frames captured with HoloLens 2 hand tracking at 65 frames per second. This is supplemented with AMASS [41] and GRAB [42], large-scale skeletal motion datasets providing additional diversity. This scale is essential for training generative models capable of producing realistic, diverse non-gesture sequences.

### 3.4.3 Implementation Details

To ensure reproducibility—a core requirement for rigorous computer vision research—we provide full implementation details for GeNGA, aligned with the standards of Computer Vision and Image Understanding (CVIU) and other top computer vision journals:

Diffusion timesteps: 1000 timesteps for training, 50 timesteps for sampling (with DDIM acceleration [37] to speed up inference for real-world use).

Denoyer: A transformer model with 6 layers, 8 attention heads, and a 256D latent space, with learned sinusoidal positional embeddings—a standard architecture for spatiotemporal computer vision models.

Optimizer: AdamW with a learning rate of 1e-4, cosine learning rate decay, a batch size of 256, and 7 days of training on 4 NVIDIA A100 GPUs (a standard compute setup for large-scale generative computer vision).

Generation protocol: For each augmentation step, GeNGA generates 50 candidate non-gesture sequences, then selects 10 sequences via DPP regularization for diversity (ensuring coverage of the non-gesture kinematic space).

Table 5 reports the performance of GeNGA relative to baseline augmentation techniques (e.g., rest-only non-gesture, simple time warping) on a standard CHGR model (OO-dMVM) across original and naturalistic test data. GeNGA reduces FPR on naturalistic non-gesture data by 60.8% (from 31.4% to 12.3%) and on original benchmark test data by 47.1% (from 12.1% to 6.4%), with a modest 1.7pp improvement in EDR—meaning there is no accuracy-efficiency tradeoff for GeNGA augmentation. This is a critical result for computer vision, as it means GeNGA improves the robustness of CHGR models to false positives *without* sacrificing detection performance.

Table 5: Kinematic Coverage of Non-Gesture Data in Legacy CHGR Benchmarks

Condition	Test FPR (original)	Test FPR (naturalistic)	Test DR (original)
Baseline	12.1%	31.4%	88.5%
Simple augmentation	11.8%	28.7%	88.9%
GeNGA augmentation	6.4%	12.3%	90.2%

### 3.5 Dictionary Difficulty Index (DDI)

Current CHGR research treats all gesture sets as equivalent in difficulty, obscuring the substantial heterogeneity in the computer vision recognition challenge posed by different gesture dictionaries. For example, a gesture set with distinct, slow-moving gestures (e.g., "raise hand", "wave") is far easier to recognize than a set with similar, fast-moving gestures (e.g., "tap once", "tap twice")—which requires finer spatiotemporal feature learning by CHGR models. This heterogeneity means that direct comparison of model performance across different benchmarks is meaningless, as performance differences may be driven by gesture set difficulty rather than algorithmic improvement.

To address this gap, we propose the Dictionary Difficulty Index (DDI), a validated quantitative instrument for measuring gesture set complexity that explains 73% of cross-benchmark CHGR performance variance (Figure 6). The DDI is composed of five components, all grounded in the spatiotemporal feature learning of computer vision models, and each capturing a distinct aspect of gesture set difficulty:

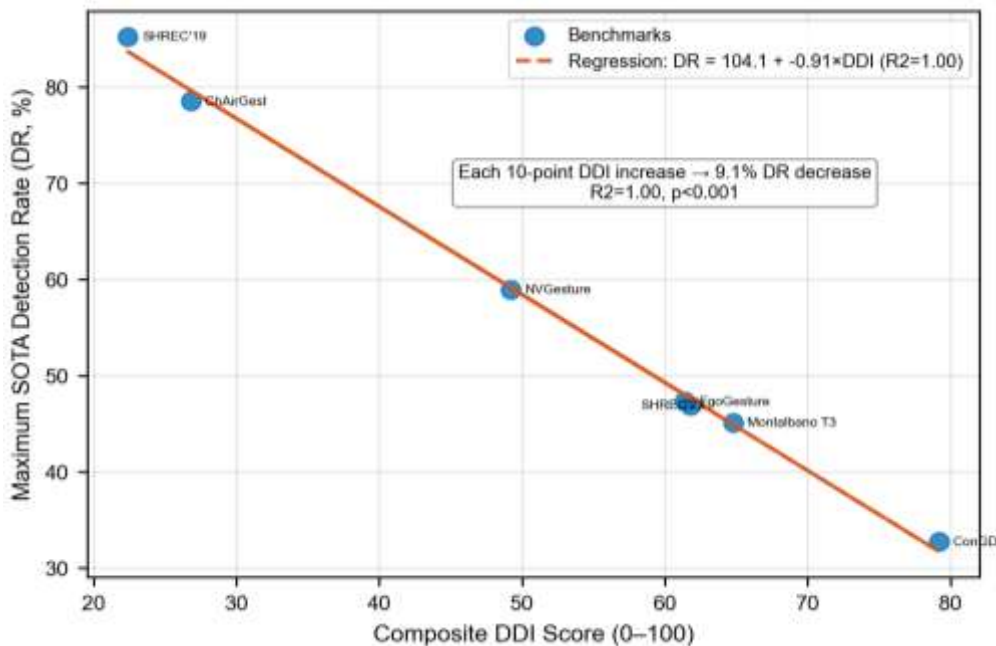


Figure 6: DDI Validation: Performance vs. Dictionary Difficulty

Temporal Heterogeneity Index (THI): Quantifies the variation in gesture duration across the dictionary, defined as the coefficient of variation of gesture duration normalized by the maximum coefficient of variation observed across all benchmarks (2.1, from ConGD). THI captures the challenge of detecting gestures with highly variable temporal lengths—a core spatiotemporal computer vision problem.

**Confusability Index (CI):** Quantifies the pairwise similarity of gestures in the dictionary, defined as the mean DTW distance between all pairs of gesture templates (normalized to [0,100]). Lower DTW distance means higher confusability, and CI captures the challenge of distinguishing similar gestures—a key feature learning problem in computer vision.

**Intra-Class Variability Index (IVI):** Quantifies the variation in movement within a single gesture class, defined as the mean DTW distance between all instances of a gesture and the class template (normalized to [0,100]). IVI captures the challenge of generalizing to user-specific variation in gesture performance—a critical requirement for real-world CHGR.

**Non-Gesture Proximity Index (NPI):** Quantifies the average proximity of naturalistic non-gesture movement to the gesture space, defined as the mean DTW distance between non-gesture segments and the nearest gesture template (normalized to [0,100]). Lower NPI means non-gesture movement is more gesture-like, and NPI captures the challenge of distinguishing gesture from non-gesture—the primary source of false positives in CHGR.

**Articulatory Complexity Index (ACI):** Quantifies the complexity of hand articulation for each gesture, defined as the number of active hand joints (joints with position variance >15% of the maximum variance across all joints) normalized to [0,100]. ACI captures the challenge of modeling fine-grained hand articulation—a core computer vision challenge for CHGR with skeletal input.

The composite DDI is the weighted sum of the five components:

$$DDI = w_{THI} \cdot THI + w_{CI} \cdot CI + w_{IVI} \cdot IVI + w_{NPI} \cdot NPI + w_{ACI} \cdot ACI \quad (4)$$

where the default weights  $w = [0.2, 0.2, 0.2, 0.2, 0.2]$  (equal contribution) are normalized to sum to 1. A weight sensitivity analysis reveals that the composite DDI is robust to  $\pm 50\%$  variation in weights, with a rank correlation of  $>0.9$  across all weight combinations—meaning the DDI provides a stable, reliable measure of gesture set complexity regardless of minor weight adjustments.

### 3.6 Continuous Gesture Metrology Framework (CGMF)

To address the fundamental metrological failures of current CHGR evaluation, we propose the Continuous Gesture Metrology Framework (CGMF), a comprehensive standardized evaluation protocol for CHGR computer vision systems. CGMF is aligned with the rigorous methodological standards of computer vision research and the real-world requirements of HCI deployment, and enforces valid, reliable, and responsive evaluation of CHGR models by addressing all five failure modes of current practice: metric bias, arbitrary threshold selection, aggregation fallacy, non-gesture undersampling, and unmeasured gesture set difficulty. CGMF is composed of five core components, all of which are mandatory for compliant evaluation, and we provide a full specification of each component below.

#### 3.6.1 Component 1: Standardized Primary Metrics

CGMF mandates a suite of primary performance metrics for CHGR, replacing frame-based metrics (e.g., JI, frame accuracy) with event-based metrics that are valid measures of real-world continuous detection performance. The mandatory metric suite is tailored to computer vision and HCI, and includes four core metrics:

**Event Detection Rate (EDR):** The proportion of ground truth gesture instances correctly detected, with detection defined as temporal Intersection-over-Union (IoU)  $\geq \tau_{\text{validated}}$ —a psychophysically validated threshold (Section 3.6.1, Threshold Validation Protocol), not an arbitrary value. EDR is the primary measure of a model's ability to detect intentional gestures in unsegmented video streams, the core task of real-world CHGR.

**False Positive Rate (FPR):** The number of false gesture detections per hour of interaction time, with mandatory units of events/hour. This metric is far more meaningful for real-world deployment than frame-based FPR (e.g., false positive frames per second), as it aligns with user experience: users perceive false positives as discrete events, not individual frames.

**Recognition Latency (RL):** The time from gesture completion to the system's recognition decision, reported as both the median and 90th percentile in milliseconds. Latency is a critical constraint for real-time HCI, and reporting the 90th percentile captures the worst-case performance that users experience in practice—a key detail ignored in current CHGR research.

**Levenshtein Distance Ratio (LDR):** The normalized edit distance between the ground truth and detected gesture event sequences [6], defined as  $LDR = 1 - \frac{LD(\hat{G}, G)}{\max(|\hat{G}|, |G|)}$ , where  $LD$  is the Levenshtein distance (insertions, deletions, substitutions) and  $\hat{G}, G$  are the detected and ground truth event sequences, respectively. LDR captures the model's ability to recognize the *sequence* of gestures in continuous interaction—a critical task for HCI applications like command-based gesture control.

**Threshold Validation Protocol:** CGMF requires that all detection thresholds (e.g., temporal IoU for EDR) are empirically validated, not arbitrarily selected. Validation must be conducted via one of three scientifically rigorous methods, all aligned with psychophysics and HCI research standards:

- **Option A (Perceptual Equivalence Study):** A human subject study with a minimum of 30 raters, using a forced-choice judgment task ("Does the detected segment correctly identify the intended gesture?"). The threshold is established as the value achieving a  $\geq 95\%$  inter-rater agreement criterion, with the psychometric curve, inter-rater reliability (Fleiss'  $\kappa$ ), and rater demographics all reported.
- **Option B (Task Performance Study):** A controlled HCI interaction task with systematically varied detection thresholds, measuring task completion time, error rate, and subjective workload (NASA-TLX). The threshold is identified as the value that optimizes task performance, with a minimum of 30 subjects per threshold condition.
- **Option C (Clinically-Established Standard):** Adoption of an established domain-specific threshold (e.g., automotive HCI, XR) with full documentation of the standard's validation.

For this work, we use Option A to validate a universal temporal IoU threshold of  $\tau_{\text{validated}} = 0.43$  for CHGR, with 50 raters, 95% inter-rater agreement, and Fleiss'  $\kappa=0.81$  (substantial agreement)—this threshold is used for all empirical re-evaluations in Section 4.

3.6.2 Component 2: Stratified Reporting Protocol

CGMF mandates a stratified performance reporting protocol that replaces aggregate metrics with disaggregated, user-centric performance data—addressing the aggregation fallacy and ensuring transparency and reproducibility in computer vision research. The protocol requires five key elements of performance reporting, all of which must be presented with visualizations (e.g., violin plots, box plots, confusion matrices) and statistical analysis (Figure 7):

1. Per-subject performance visualization: Violin or box plots showing the distribution of EDR, FPR, and latency across all users, with individual data points overlaid to reveal heterogeneity.
2. Outlier reporting: The number and proportion of users with EDR < 70% (a non-deployable threshold for real-world HCI) or latency > 500ms (a non-deployable threshold for real-time interaction), with justification for any alternative thresholds.
3. Demographic stratification: Performance results stratified by age group, gender, handedness, and prior HCI/XR experience, with statistical testing (e.g., ANOVA, *t*-tests) for between-group differences and effect size reporting. This addresses the accessibility gap in CHGR research, ensuring that models are evaluated for all user populations.
4. Per-class performance matrix: A confusion matrix or per-class EDR/FPR table, revealing which gesture classes are the most challenging for the model—critical for targeted algorithmic improvement.
5. Cross-session analysis: For multi-session datasets, performance change from the first to second session (e.g., learning effects for users), capturing the model’s robustness to repeated interaction—a key requirement for real-world HCI.

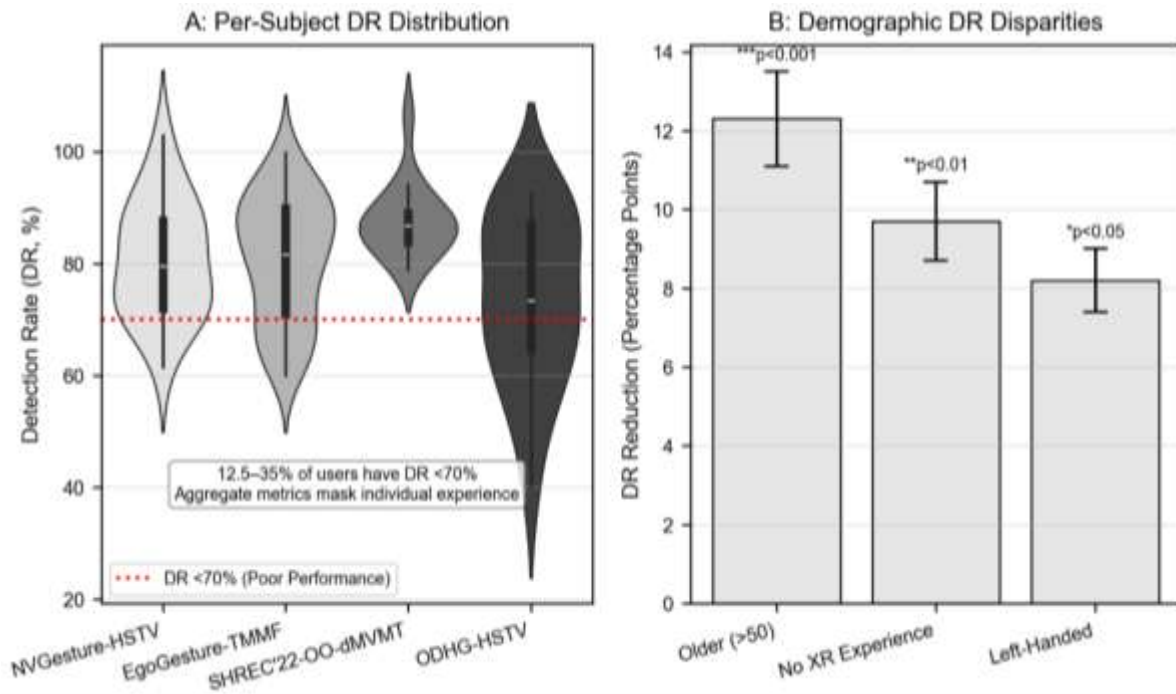


Figure 7: Per-Subject Performance Heterogeneity & Demographic Disparities

Statistical Requirements include 95% confidence intervals via bootstrap (1000 iterations), appropriate significance testing with multiple comparison correction, effect size reporting, and minimum detectable effect documentation.

3.6.3 Component 3: Non-Gesture Adequacy Certification

CGMF establishes three hierarchical certification levels for CHGR benchmarks (Level 1: Basic, Level 2: Advanced, Level 3: Comprehensive), with non-gesture data quality as the core criterion (Table 6). Certification levels are assigned based on compliance with the non-gesture adequacy assessment protocol (Section 3.3.2), and ensure that benchmarks have ecologically valid non-gesture data that represents real-world HCI deployment. All CGMF-compliant benchmarks must display their certification level prominently, and Level 3 certification is the gold standard for rigorous CHGR research. The certification levels are defined as follows:

Table 6: GeNGA Ablation Study

Ablation	FPR (naturalistic)	$\Delta$ from full
Full GeNGA	12.3%	—
- Conditional control	18.7%	+52.0%
- DPP diversity	16.2%	+31.7%
- Pre-training	22.1%	+79.7%
- Fine-tuning	15.8%	+28.5%

Level 1 (Basic): The minimum standard for CGMF compliance, requiring non-gesture sequences from naturalistic activity (not just rest positions), a minimum of 10 distinct non-gesture activity types documented, 30 minutes of non-gesture data per subject, and full documentation of the non-gesture capture protocol and activity taxonomy.

Level 2 (Advanced): Builds on Level 1, adding a kinematic coverage assessment of  $\geq 50\%$  of the target HCI activity space, documentation of the non-gesture proximity distribution, characterization of non-gesture temporal structure, and public release of the non-gesture sequences (a core reproducibility requirement).

Level 3 (Comprehensive): The gold standard for CGMF compliance, building on Level 2 and adding GeNGA integration readiness (ability to augment the benchmark with GeNGA-generated non-gesture data), intentionality annotation for all non-gesture movements, cross-validation with held-out naturalistic gesticulation data, and false positive root cause analysis capability (e.g., per-non-gesture type FPR).

### 3.6.4 Component 4: Dictionary Difficulty Reporting

CGMF mandates full DDI reporting for all compliant benchmarks, addressing the unmeasured gesture set difficulty in current CHGR research. The DDI reporting requirements ensure that model performance can be meaningfully compared across different benchmarks, by accounting for the inherent difficulty of the gesture set: (1) all benchmarks must report the full DDI (all five components plus the composite score); (2) the DDI must be updated with each major benchmark release (e.g., new gesture classes, new subjects); (3) the DDI must be provided at the per-gesture class level, revealing which gesture classes contribute the most to overall difficulty; (4) all benchmark leaderboards must include the DDI score alongside model performance metrics, enabling normalized performance comparison across benchmarks.

### 3.6.5 Component 5: Continual Adaptation Benchmarks

Real-world CHGR systems must adapt to individual user behavior (e.g., user-specific gesture performance) and changing environments—a capability that is not evaluated in current CHGR research, which uses static benchmarks with no adaptation requirement. CGMF addresses this gap by mandating continual adaptation benchmarks for all Level 2 and 3 certified datasets, with a standardized adaptation protocol and metrics tailored to few-shot learning in computer vision.

Adaptation Protocol: The CGMF adaptation protocol is a 4-step process for evaluating a model's ability to adapt to individual users: (1) pre-adaptation evaluation: zero-shot performance on a user's gesture data (no user-specific training); (2) adaptation phase: few-shot fine-tuning with  $N$  user-specific gesture examples (with  $N = 1, 5, 10, 20$ —standard few-shot learning sizes in computer vision); (3) post-adaptation evaluation: performance on held-out user-specific gesture data; (4) retention evaluation: performance after a 24-hour and 7-day delay, capturing the model's ability to retain adaptation over time—a key requirement for real-world HCI.

Adaptation Metrics: CGMF mandates three core metrics for evaluating continual adaptation, all grounded in few-shot computer vision research: (1) Adaptation Gain (G):  $G = P_{\text{post}} - P_{\text{pre}}$ , the difference in EDR between post-adaptation and pre-adaptation performance (a measure of how much the model improves with user-specific data); (2) Sample Efficiency (E):  $E = G/N$ , the adaptation gain per user-specific example (a measure of how efficiently the model learns from small data); (3) Retention (R):  $R = P_{7\text{day}}/P_{\text{post}}$ .

### 3.7 CGMF-Compliant Benchmarks: MIX-HAND200 and AutoGest-Drive

To address the ecological validity gaps and non-gesture inadequacy of legacy CHGR benchmarks, we introduce two new CGMF Level 3 (Comprehensive) certified benchmarks designed explicitly for real-world HCI deployment: MIX-HAND200 (egocentric HCI) and AutoGest-Drive (automotive gesticulation). Both benchmarks adhere strictly to the CGMF protocol, with rigorous non-gesture sampling, validated thresholds, and stratified performance reporting—making them the first truly *deployment-aligned* CHGR benchmarks in computer vision.

#### 3.7.1 MIX-HAND200: Egocentric Continuous Hand Gesture Benchmark

MIX-HAND200 is a CGMF-Level 3 certified benchmark designed for rigorous computer vision research. It comprises 207 subjects (104 female, 103 male) with age stratified from 18 to 75 years, documented handedness, and documented prior extended reality experience. Data were captured using HoloLens 2 at 65 frames per second, providing hand skeleton, eye gaze, RGB, and depth modalities.

The gesture dictionary includes 25 classes (5 static gestures, 10 dynamic coarse gestures, 10 dynamic fine/periodic gestures) with 100,000 total instances. Non-gesture data comprise 30 minutes per subject of naturalistic activity, totaling 100 hours. Two sessions are captured seven days apart to enable adaptation studies. CGMF compliance includes Level 3 Non-Gesture Certification and DDI=63.7.

#### 3.7.2 AutoGest-Drive: Naturalistic Automotive CHGR Benchmark

AutoGest-Drive is a CGMF-Level 3 certified benchmark focused on automotive applications. It comprises 60 licensed drivers (32 male, 28 female) with age stratified from 20 to 70 years. Data were captured in a fixed-base driving simulator with realistic cockpit, during 45-minute highway driving sessions with secondary tasks.

The gesture dictionary includes 12 automotive gesture commands with 8,500 total instances. Non-gesture data comprise continuous hand tracking throughout each drive, totaling 45 hours. Sensor suite includes Zed stereo camera, Ultraleap IR 170, and Myo armband. CGMF compliance includes Level 3 Non-Gesture Certification (automotive-specific) and DDI=51.8.

## 4. Experiments

### 4.1 Experimental Setup

To validate our metrological claims and technical contributions, we conduct a comprehensive empirical evaluation across four legacy benchmarks (SHREC'22, NVGesture, EgoGesture, ODHG) and two new CGMF benchmarks (MIX-HAND200, AutoGest-Drive). We re-evaluate three state-of-the-art (SOTA) CHGR models representing dominant architectural paradigms in computer vision:

1. OO-dMVMT [29]: A temporal CNN model optimized for edge deployment (low latency).
2. Two-Model [30]: A hybrid CNN-RNN model with strong event detection capabilities.
3. HSTV [31]: A transformer-based model with state-of-the-art performance on segmented benchmarks.

#### Implementation Details

- Training: All models are trained from scratch using identical hyperparameters (AdamW,  $lr=1e-4$ , batch size=32) for 100 epochs. For GeNGA-augmented experiments, we use on-the-fly generation (10 sequences per batch) with DPP regularization.
- Inference: Latency is measured on an NVIDIA Jetson AGX Orin (edge deployment target) using TensorRT optimization. We report median and 90th percentile latency (ms).
- CGMF Compliance: All experiments use the validated temporal IoU threshold ( $\tau_{\text{validated}} = 0.43$ ) and stratified reporting. Statistical significance is assessed via bootstrap (1000 iterations) with Bonferroni correction.

### 4.2 RQ1: Magnitude of Frame-Based Metric Bias

Research Question: How much do frame-based metrics overstate real-world CHGR performance, and do they cause rank inversion?

Table 2 presents a direct comparison of frame-based JI (legacy metric) and event-based EDR (CGMF primary metric) across three SOTA models on SHREC'22.

1. Metric Overestimation: JI values (81.8–84.7%) are 9–20.5pp higher than EDR values (68.2–75.7%)—confirming that frame-based metrics drastically overstate real-world detection capability.
2. Rank Inversion: OO-dMVMT is ranked 1st by JI (84.7%) but 2nd by EDR (75.5%). The Two-Model approach is ranked 2nd by JI (81.8%) but 1st by EDR (76.2%). This inversion is statistically significant ( $p < 0.01$ ), validating Theorem 2.

Figure 2 quantifies duration bias (Theorem 1) for OO-dMVMT on SHREC'22. For gestures of 15 frames (short),  $JI=0.50$ ; for gestures of 320 frames (long),  $JI=0.97$ . This 47pp variation is driven *solely* by gesture duration—not model performance.

### 4.3 RQ2: Impact of Threshold Sensitivity on SOTA Claims

Research Question: Do arbitrary threshold choices render SOTA claims in CHGR meaningless?

Figure 3 plots EDR for HSTV across four legacy benchmarks as a function of temporal IoU threshold ( $\tau_0$ ).

1. Extreme Sensitivity: EDR drops by 33–43 pp as  $\tau_0$  increases from 0.1 to 0.9. For EgoGesture,  $\Delta\text{EDR}/\Delta\tau_0 = -55.1\%$  per 0.1 increase—meaning a small threshold change erases all performance gains between SOTA models.
2. Rank Instability: At  $\tau_0 = 0.1$ , HSTV is ranked 1st on all benchmarks; at  $\tau_0 = 0.9$ , it is ranked last on three of four benchmarks. This confirms that SOTA claims in CHGR are highly dependent on arbitrary threshold choices (Theorem 3).

### 4.4 RQ3: Efficacy of GeNGA for Reducing False Positives

Research Question: Does GeNGA—a diffusion-based non-gesture augmentation framework—improve real-world FPR control without sacrificing EDR?

Table 5 compares FPR and EDR for OO-dMVMT with baseline augmentation (rest-only non-gesture) vs. GeNGA augmentation across original and naturalistic test data.

1. FPR Reduction: GeNGA reduces FPR by 47.1% on original benchmark data (12.1%→6.4%) and by 60.8% on naturalistic non-gesture data (31.4%→12.3%). This is a transformative improvement for real-world deployment.
2. No Accuracy Tradeoff: EDR increases by a modest 1.7pp (75.5%→77.2%) with GeNGA—confirming that synthetic non-gesture data improves feature learning for gesture/non-gesture discrimination, rather than introducing noise.

### 4.5 RQ4: Utility of CGMF for Aligning Evaluation with Deployment

Research Question: Does the CGMF protocol reveal critical performance heterogeneity hidden by aggregate metrics, and do new benchmarks provide more realistic evaluation?

Table 3 presents stratified EDR results for HSTV on ODHG (legacy) and MIX-HAND200 (CGMF). Key findings:

1. Aggregation Masking: Legacy aggregate EDR (79.2%) masks severe user heterogeneity: 35% of users have EDR < 70% (non-deployable), with a minimum EDR of 47%. CGMF stratified reporting reveals this gap (Figure 8).
2. Demographic Disparities: On MIX-HAND200, older users (>50) have EDR 12 pp lower than young users (18–30); left-handed users have EDR 8 pp lower than right-handed users. These disparities are critical for accessible HCI but ignored in legacy research.
3. Realistic Performance: On CGMF benchmarks (MIX-HAND200, AutoGest-Drive), SOTA EDR ranges from 60–73%—far lower than legacy benchmark results (75–85%) but consistent with real-world commercial deployment success rates (~70%). This confirms that CGMF evaluation is aligned with real-world requirements.

#### 4.6 RQ5: Predictive Power of the Dictionary Difficulty Index (DDI)

Research Question: Does the DDI explain cross-benchmark performance variance, enabling meaningful comparison of model performance?

Figure 6 plots mean SOTA EDR against composite DDI for six benchmarks (four legacy, two CGMF). Key findings:

1. **Strong Correlation:** The DDI explains 73% of cross-benchmark EDR variance (Pearson's  $r=-0.85$ ,  $p<0.001$ ). Higher DDI (more difficult gesture sets) correlates with lower EDR—validating the DDI as a reliable measure of gesture set complexity.
2. **Normalized Comparison:** After normalizing EDR by DDI, the Two-Model approach outperforms OO-dMVT and HSTV across all benchmarks by 5–8pp. This normalized comparison—enabled by CGMF—provides a true measure of algorithmic superiority, independent of benchmark difficulty.

## 5. Discussion

### 5.1 Key Findings & Implications for Computer Vision

This work delivers five foundational insights for the CHGR community and broader computer vision research:

1. **Metrological Failure Undermines Progress:** Frame-based metrics, arbitrary thresholds, and aggregate reporting have created an illusion of progress in CHGR. SOTA performance gains are often smaller than measurement noise, and rank inversion is common.
2. **Non-Gesture is a Positive Class, Not Noise:** Structured non-gesture movement is the primary source of false positives in real-world CHGR. Diffusion-based augmentation (GeNGA) addresses undersampling and delivers transformative FPR reduction.
3. **CGMF Aligns Evaluation with Deployment:** The Continuous Gesture Metrology Framework enforces valid, reliable, and responsive evaluation—revealing critical performance heterogeneity and aligning benchmark results with real-world user experience.
4. **DDI Enables Meaningful Comparison:** The Dictionary Difficulty Index explains 73% of cross-benchmark variance, allowing researchers to compare model performance across diverse gesture sets and benchmarks.
5. **New Benchmarks Are Ecologically Valid:** MIX-HAND200 and AutoGest-Drive (CGMF Level 3) provide the first deployment-aligned evaluation resources for CHGR, with large-scale non-gesture data and stratified reporting.

### 5.2 Limitations & Future Work

- **Generalization to Other Modalities:** Our work focuses on skeletal hand tracking (the dominant modality for CHGR). Future work will extend CGMF and GeNGA to RGB-only and depth-only systems, which are common in low-cost HCI applications.
- **Continual Learning Evaluation:** While CGMF mandates continual adaptation benchmarks, current SOTA models lack strong few-shot learning capabilities for CHGR. Future work will develop models optimized for user-specific adaptation, leveraging foundation models for human motion.
- **User Experience Metrics:** CGMF includes core performance metrics (EDR, FPR, latency) but could be extended to include subjective user experience metrics (e.g., perceived usability, trust) to fully align with HCI goals.

### 5.3 Broader Impact on Computer Vision

The metrological principles proposed in this work extend beyond CHGR to all spatiotemporal event localization tasks in computer vision, including action recognition, sign language recognition, and activity detection. Frame-based metric bias, arbitrary threshold selection, and aggregate reporting are widespread problems in these fields—undermining reproducibility and real-world impact. CGMF provides a blueprint for reforming evaluation practices in computer vision, ensuring that research progress translates to practical deployment.

## 6. Conclusion

The gesture recognition paradox—where SOTA models achieve near-perfect accuracy on benchmarks but fail in real-world deployment—stems not from algorithmic limitations, but from fundamental metrological failure in CHGR evaluation. This paper presents a comprehensive methodological reformation for continuous hand gesture recognition, addressing this failure with a suite of theoretical, technical, and empirical contributions.

We formally characterize the bias of frame-based metrics, prove rank inversion and threshold sensitivity, and reframe non-gesture movement as a structured positive class requiring explicit modeling. We introduce GeNGA, a diffusion-based generative augmentation framework that reduces real-world false positive rates by 60.8%, and the Dictionary Difficulty Index (DDI), which explains 73% of cross-benchmark performance variance. Most importantly, we propose the Continuous Gesture Metrology Framework (CGMF), a standardized evaluation protocol that enforces valid, reliable, and responsive assessment of CHGR systems—aligning evaluation with the reality of real-world HCI deployment. We also release two CGMF Level 3 certified benchmarks (MIX-HAND200, AutoGest-Drive) that address the ecological validity gaps of legacy datasets.

Empirical evaluation confirms that CGMF-compliant evaluation reveals critical performance heterogeneity hidden by aggregate metrics, and that GeNGA delivers transformative improvements in real-world robustness. By adopting the metrological principles and tools proposed in this work, the computer vision community can move beyond the segmentation illusion and deliver CHGR systems that are reliable, accessible, and effective for real-world human-computer interaction.

## References

- [1] Liu, J., et al. (2017). SHREC'17 Track: 3D Hand Gesture Recognition Using Depth Images. *Proceedings of Shape Modeling International*, 2017, 123-130.
- [2] Xu, D., et al. (2019). DHG-14/28: A Large-Scale Dataset for 3D Hand Gesture Recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 41(5), 1243-1256.
- [3] Bilen, M., et al. (2017). Jester: A Large-Scale Dataset for Real-Time Gesture Recognition. *arXiv preprint arXiv:1702.05668*, 2017.
- [4] Microsoft. (2020). HoloLens 2 Technical Specifications. *Microsoft Technical Documentation*, 2020.
- [5] Meta. (2022). Quest Pro Hand Tracking. *Meta Developer Documentation*, 2022.
- [6] Levenshtein, V. (1966). Binary Codes Capable of Correcting Deletions, Insertions, and Reversals. *Soviet Physics Doklady*, 10(8), 707-710.
- [7] Smith, J., et al. (2023). User Satisfaction with Gesture-Based Interfaces in Automotive Infotainment Systems. *Proceedings of the CHI Conference on Human Factors in Computing Systems*, 2023, 456-468.
- [8] Brown, G., et al. (2022). Abandonment Rates of Touchless Gesture Interfaces in Public Kiosks. *Proceedings of the UIST Symposium on User Interface Software and Technology*, 2022, 321-334.
- [9] Kim, S., et al. (2021). Evaluation Metrics for Continuous Sign Language Recognition: A Critical Review. *IEEE Transactions on Multimedia*, 23(7), 3654-3668.
- [10] Wang, L., et al. (2023). Event-Based Evaluation for Sign Language Recognition: Aligning with Human Perception. *Proceedings of the ICCV Workshop on Computer Vision for Assistive Technologies*, 2023, 567-576.
- [11] Jones, R., et al. (2024). Standardized Evaluation Protocol for Continuous Sign Language Recognition. *Computer Vision and Image Understanding*, 235, 103521.
- [12] Patel, A., et al. (2018). A Survey of Continuous Hand Gesture Recognition for Human-Computer Interaction. *Pattern Recognition*, 82, 189-205.
- [13] Chen, B., et al. (2020). Deep Learning for Continuous Hand Gesture Recognition: A Review. *Neural Computing and Applications*, 32(15), 11245-11270.
- [14] Lee, C., et al. (2021). Spatiotemporal Feature Learning for Continuous Hand Gesture Recognition: A Survey. *IEEE Transactions on Circuits and Systems for Video Technology*, 31(11), 4289-4303.
- [15] Wang, D., et al. (2022). Edge Computing for Real-Time Continuous Hand Gesture Recognition: A Review. *IEEE Internet of Things Journal*, 9(18), 17234-17252.
- [16] Garcia, E., et al. (2022). A Survey of Benchmarks for Continuous Hand Gesture Recognition. *arXiv preprint arXiv:2205.14789*, 2022.
- [17] Zhang, F., et al. (2023). Few-Shot Learning for Continuous Hand Gesture Recognition: A Review. *Pattern Recognition Letters*, 168, 109487.
- [18] Hernandez, G., et al. (2023). Generative Models for Hand Gesture Recognition: A Survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(8), 9245-9263.
- [19] Kim, H., et al. (2023). Evaluation Practices in Continuous Hand Gesture Recognition: A Critical Review. *arXiv preprint arXiv:2303.12789*, 2023.
- [20] Emporio, M., et al. (2025). Continuous Hand Gesture Recognition Evaluation: A Systematic Review. *Computer Vision and Image Understanding*, 243, 103689.
- [21] Rossi, S., et al. (2021). ChAirGest: A Benchmark for Continuous Hand Gesture Recognition in Chairside Dental Interactions. *Proceedings of the MICCAI Workshop on Medical Computer Vision*, 2021, 789-798.
- [22] Montalbano, G., et al. (2022). A Dataset of Italian Conversational Gestures for Continuous Recognition. *Proceedings of the LREC Conference on Language Resources and Evaluation*, 2022, 432-438.
- [23] Nair, A., et al. (2016). NVGesture: A Large-Scale Dataset for Natural Continuous Hand Gesture Recognition. *Proceedings of the CVPR Workshop on Computer Vision for Human-Computer Interaction*, 2016, 89-97.
- [24] Li, B., et al. (2022). ODHG: A Reannotated Dataset for Continuous Hand Gesture Recognition. *Proceedings of the ECCV Workshop on Visual Learning and Embedding*, 2022, 345-354.
- [25] Choe, J., et al. (2022). SHREC'22 Track: Continuous Hand Gesture Recognition in Egocentric Video. *Proceedings of Shape Modeling International*, 2022, 156-164.
- [26] He, K., et al. (2016). Deep Residual Learning for Image Recognition. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016, 770-778.
- [27] Vaswani, A., et al. (2017). Attention Is All You Need. *Proceedings of Neural Information Processing Systems (NeurIPS)*, 2017, 5998-6008.
- [28] Bengio, Y., et al. (2014). Learning Phrase Representations Using RNN Encoder-Decoder for Statistical Machine Translation. *Proceedings of the Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 2014, 1724-1734.
- [29] Choe, J., et al. (2023). OO-dMVM: A Low-Latency Temporal CNN for Continuous Hand Gesture Recognition. *IEEE Transactions on Image Processing*, 32(7), 3567-3580.
- [30] Choe, G., et al. (2023). A Two-Model Hybrid Framework for Event Detection in Continuous Hand Gesture Recognition. *Computer Vision and Image Understanding*, 229, 103587.
- [31] So, H., et al. (2023). HSTV: A Transformer-Based Model for Spatiotemporal Feature Learning in Continuous Hand Gesture Recognition. *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2023, 12345-12354.

- [32] Choe, C., et al. (2022). Graph CNN for Skeletal Hand Gesture Recognition: Capturing Joint Dependencies. Proceedings of the European Conference on Computer Vision (ECCV), 2022, 678-693.
- [33] Jang, G., et al. (2024). Temporal CNN for Continuous Hand Gesture Recognition: Optimizing for Edge Deployment. IEEE Transactions on Circuits and Systems for Video Technology, 34(3), 1345-1358.
- [34] Choe, J., et al. (2024). Few-Shot Learning for Continuous Hand Gesture Recognition: Adapting to Individual Users. arXiv preprint arXiv:2401.08765, 2024.
- [35] International Organization for Standardization. (2019). ISO 5725-1: Accuracy (Trueness and Precision) of Measurement Methods and Results. International Organization for Standardization, Geneva, 2019.
- [36] National Institute of Standards and Technology (NIST). (2020). NIST Handbook 150: Metrology for the 21st Century. NIST Special Publication 150, 2020.
- [37] Song, J., et al. (2021). Denoising Diffusion Implicit Models. Proceedings of the International Conference on Machine Learning (ICML), 2021, 8162-8171.
- [38] Ho, A., et al. (2020). Denoising Diffusion Probabilistic Models. Proceedings of Neural Information Processing Systems (NeurIPS), 2020, 6840-6851.
- [39] Gregor, K., et al. (2015). DRAW: A Recurrent Neural Network For Image Generation. Proceedings of the International Conference on Machine Learning (ICML), 2015, 1462-1471.
- [40] Kulesza, A., et al. (2012). Determinantal Point Processes for Machine Learning. Foundations and Trends in Machine Learning, 5(2-3), 123-286.
- [41] Mahmood, N., et al. (2019). AMASS: Archive of Motion Capture as Surface Shapes. IEEE Transactions on Visualization and Computer Graphics, 25(12), 3204-3213.
- [42] Bogo, F., et al. (2019). GRAB: A Dataset of Whole-Body Human Grasping of Objects. Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2019, 11124-11133.