# Zero-Shot Open-Vocabulary Human Motion Grounding with Test-Time Training

**Yunjiao Zhou, Xinyan Chen, Junlang Qian, Lihua Xie, Jianfei Yang** *

Nanyang Technological University, Singapore
{yunjiao001, chen1909, junlang001}@e.ntu.edu.sg, {elhxie, jianfei.yang}@ntu.edu.sg

## Abstract

Understanding complex human activities demands the ability to decompose motion into fine-grained, semantic-aligned sub-actions. This motion grounding process is crucial for behavior analysis, embodied AI and virtual reality. Yet, most existing methods rely on dense supervision with predefined action classes, which are infeasible in open-vocabulary, real-world settings. In this paper, we propose ZOMG, a zero-shot, open-vocabulary framework that segments motion sequences into semantically meaningful sub-actions without requiring any annotations or fine-tuning. Technically, ZOMG integrates (1) language semantic partition, which leverages large language models to decompose instructions into ordered sub-action units, and (2) soft masking optimization, which learns instance-specific temporal masks to focus on frames critical to sub-actions, while maintaining intra-segment continuity and enforcing inter-segment separation, all without altering the pretrained encoder. Experiments on three motion-language datasets demonstrate state-of-the-art effectiveness and efficiency of motion grounding performance, outperforming prior methods by +8.7% mAP on HumanML3D benchmark. Meanwhile, significant improvements also exist in downstream retrieval, establishing a new paradigm for annotation-free motion understanding.

**Code** — https://github.com/pridy999/ZOMG

## Introduction

Understanding human motion at a fine-grained level is vital for tasks such as behavior analysis, embodied AI, and virtual reality (Zhang et al. 2023; Zhou, Wan, and Wang 2024; Zhou et al. 2023a; Yang et al. 2022b). However, real-world motion is temporally unstructured and composed of overlapping sub-actions without explicit boundaries. This absence of temporal modularity hinders downstream applications such as motion retrieval (Xue et al. 2025; Zhou et al. 2024) and generation (Zeng et al. 2025), which rely on semantically meaningful motion units for reasoning and interaction. In such cases, motion grounding, which decomposes continuous motion streams into semantically coherent units, plays a pivotal role in enabling structured motion representations and flexible motion interaction.
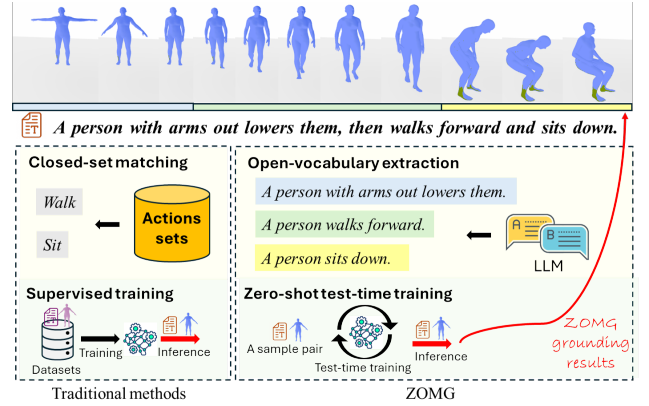
---

*Corresponding author.

Figure 1: Motion grounding illustration of ZOMG.

Existing motion grounding methods (Chen, Tsai, and Yang 2021; Yang et al. 2022a) typically rely on closed-set assumptions (Zhou et al. 2023b; Radford et al. 2021), training models to align motion sequences with labels from a fixed action vocabulary (e.g., "walk", "wave"). However, real-world motion is inherently open-ended and compositional, with free-form expressions describing subtle, overlapping sub-actions (e.g., "a person sits down while waving") that cannot be captured by predefined atomic labels. This rigid setup fundamentally limits scalability and generalization, highlighting the need for open-vocabulary grounding that aligns motion with natural language. In addition, the scarcity of fine-grained annotations and the diversity of real-world expressions motivate a zero-shot setting (Punnakkal et al. 2021; Qian et al. 2025), where models are required to generalize to novel queries without task-specific supervision. Therefore, our research focus is *to achieve reliable and efficient motion grounding under a zero-shot open-vocabulary setting.*

However, realizing zero-shot open-vocabulary motion grounding poses several challenges. First, open-world descriptions are free-form and structurally ambiguous, often lacking explicit cues regarding the number of sub-actions and their complex temporal interactions (Liu et al. 2022). For instance, phrases like "walks confidently and waving" mix concurrent and sequential elements, making it difficult for rule-based models (Wu et al. 2022) to extract correct segments. Second, while existing motion-language mod-

els (Tevet et al. 2022) capture global semantics, grounding requires finer temporal resolution to detect subtle transitions between sub-actions. This necessitates frame-level reasoning, which current sequence-level encoders are not inherently equipped to handle. Third, even when frame-level features are available, motion exhibits strong temporal continuity and entanglement. Unlike static images, individual frames often lack standalone semantic identity, and sub-actions emerge only through inter-frame dynamics. Without modeling these dependencies, frame-wise analysis alone is insufficient for accurate grounding.

Given these challenges, it is necessary to obtain finer-grained representations for both text and motion. On the language side, we compensate for the absence of segment-level annotations by using large language models (LLMs) (Zhao et al. 2023) to decompose free-form descriptions into ordered sub-action units. These unit-level queries provide semantic guidance for open-vocabulary matching. On the motion side, directly supervising frame-level alignment is impractical due to annotation scarcity. Pretrained motion-language models, despite being trained only at the sequence level, implicitly encode the inter-frame transitions necessary for grounding (Wang et al. 2023a). This suggests that exploiting such temporal structure already embedded in the representation space is promising for fine-grained motion grounding. To make this structure accessible without altering model parameters, we design a test-time training approach. It introduces a small set of variables for each input, which are optimized to selectively attend to sub-action-relevant frames. This adaptation refines the attention behavior of pretrained model (Lee, Lee, and Kang 2019) on a per-instance basis, allowing precise localization of motion segments while maintaining generalization in zero-shot settings.

In this work, we propose ZOMG, a test-time grounding framework that extracts fine-grained motion structure from pretrained motion-language models without annotations. ZOMG bypasses the need for fine-tuning by introducing a lightweight test-time training stage that transforms sequence-level representations into temporally grounded sub-actions. It consists of two modules: (1) Language Semantic Partition (LSP) uses an LLM to decompose free-form textual descriptions into semantically coherent and temporally ordered sub-action queries, serving as anchors for subsequent grounding; (2) Soft Masking Optimization (SMO) aligns each query with motion by optimizing frame-wise soft masks, guided by semantic alignment and structural regularizations to ensure segment separability and continuity. Experiments across large-scale datasets demonstrate that ZOMG improves grounding accuracy by up to +8.7% mAP on HumanML3D and significantly boosts downstream motion-text retrieval. Despite these gains, it remains highly efficient, requiring only 0.5K optimization parameters and achieving over 3× inference speedup compared to existing TTT methods, enabling practical annotation-free deployment.

The contributions are summarized as follows:

- We introduce ZOMG, the first framework that enables annotation-free and open-vocabulary motion grounding, achieving +8.7% mAP improvement on HumanML3D.

- We propose a novel test-time training scheme combining LLM-guided decomposition and soft masking optimization, uncovering the fine-grained temporal structure from pretrained motion-language models.

- Extensive experiments show that ZOMG achieves state-of-the-art grounding accuracy and significantly improves motion-text retrieval. Despite these gains, it remains highly efficient at test time, supporting its practical deployment in real-world settings.

## Related Work

### Temporal Grounding

Temporal grounding seeks to localize video segments that correspond to free-form textual descriptions (Chao et al. 2018; Zhao et al. 2017; Yan et al. 2023; Nguyen et al. 2025). Early approaches like TALL (Gao et al. 2017) and 2D-TAN (Zhang et al. 2020) formulate this task as moment retrieval via cross-modal attention and sliding-window proposals, trained under dense annotations and closed-set vocabularies. Recent work has explored open-vocabulary and zero-shot settings by leveraging pretrained vision-language models. For example, T3AL (Liberatori et al. 2024) uses CLIP-based retrieval to match language queries with video frames, while STALE (Nag et al. 2022) introduces a one-stage model with parallel streams for localization and classification. X-POOL (Gorti et al. 2022) and SeViLA (Yu et al. 2023) further enhance grounding via frozen VLMs (Xu et al. 2021; Bordes et al. 2024), augmented with adapters or prompting strategies. These methods showcase the power of pretrained representations for temporal localization, but are largely limited to video domains, where large-scale vision-language models provide strong priors. In contrast, motion data lacks such pretrained encoders, making zero-shot grounding in this space substantially more challenging.

### Motion Understanding

Understanding and generating human motion from language has drawn increasing interest for applications in animation, robotics, and embodied AI. Sequence-level models like TEMOS (Petrovich, Black, and Varol 2022), TMR (Petrovich, Black, and Varol 2023), and TEACH (Athanasiou et al. 2022) learn joint motion-text representations via contrastive or generative training. Others, including MotionCLIP (Tevet et al. 2022) and Motion-X (Lin et al. 2023), incorporate CLIP priors to encode human motion, but operate at the sequence level, limiting their capacity for segment-level reasoning. Recent studies have begun to explore finer temporal structure through action segmentation (Wang et al. 2023b; Kong et al. 2019; Punnakkal et al. 2021) and fine-grained motion retrieval (Li and Feng 2024; Wang, Kang, and Mu 2024; Zhang et al. 2023). However, these methods typically depend on predefined taxonomies or task-specific labels, and seldom address frame-level grounding in open settings. In contrast, our work aims to understand fine-grained motion under open-vocabulary, zero-shot conditions. Without predefined categories, our method extracts semantically coherent motion units that can be reused across diverse tasks, enhancing both generalization and interpretability.

## Method

### Problem Definition

Zero-shot open-vocabulary human motion grounding aims to segment compositional motion sequences into semantically coherent units aligned with free-form text, without predefined action classes or temporal annotations. Formally, given a text $\mathcal{T}$ and a motion sequence $\mathcal{M} = \{m_1, \ldots, m_L\}$ of $L$ frames, the goal is to segment $\mathcal{M}$ into $\{S_1, \ldots, S_k\}$, where each segment $S_i$ corresponds to a sub-action aligned with a span $\mathcal{T}_i \subseteq \mathcal{T}$. Our framework, shown in Fig. 2, includes: (1) Motion-Language Pretraining, which establishes a fundamental understanding of global motion-text semantic alignment, and (2) Test-Time Grounding, which adapts the pretrained model on a per-instance basis to segment motion sequences in alignment with the input text.

### Motion-Language Pretraining

Motion-language pretraining focuses on sequence-level alignment, implicitly capturing inter-frame temporal dynamics to provide cues for downstream segmentation even without segment-level supervision. Given paired motion-text samples $(\mathcal{M}, \mathcal{T})$, the motion encoder $E_M$ maps $\mathcal{M}$ to per-frame features $F = \{f_1, \ldots, f_L\}, f_i \in \mathbb{R}^d$, while the text encoder $E_T$ produces a global text embedding $t \in \mathbb{R}^d$. To retain fine-grained temporal cues, $E_M$ explicitly computes these frame-level representations before any temporal aggregation. An attention pooling module $P(\cdot)$ then computes the motion embedding $m \in \mathbb{R}^d$ via learned weights, integrating local semantics with temporal salience. The model is trained with a contrastive loss (Radford et al. 2021):

$$\mathcal{L}_c = -\log \frac{\exp(m \cdot t^+ / \tau)}{\sum_{t \in \{t^+, t^-\}} \exp(m \cdot t / \tau)}, \qquad (1)$$

where $t^+$ and $t^-$ denote positive and negative text samples, respectively. This pretraining yields general-purpose encoders that implicitly capture motion structure, laying the foundation for zero-shot fine-grained grounding.

### Test-Time Grounding

While existing approaches often rely on fine-tuning (Wortsman et al. 2022) or adapter-based training (Wang et al. 2020) to specialize pretrained models for grounding tasks, such methods require nontrivial amounts of segment-level supervision and may overfit to training distributions. This contradicts our goal of fully zero-shot, open-vocabulary motion grounding. To address this, we adopt a test-time training strategy (Sun et al. 2020; Liu et al. 2021), which enables instance-specific adaptation without modifying the pretrained model. It introduces: (1) LSP, which leverages LLMs to decompose free-form textual descriptions into temporally ordered sub-action units, providing structured semantic anchors for grounding. (2) SMO, which performs instance-specific optimization to generate frame-wise soft masks, identifying sub-action-relevant frames without supervision.

**Language Semantic Partition** Zero-shot motion grounding lacks segment-level annotations, making it difficult to localize complex instructions over time. Traditional rule-based or syntactic parsing struggles with ambiguity and compositionality in open-vocabulary descriptions (FitzGerald et al. 2018). To bridge this gap, we harness LLMs' reasoning to decompose complex descriptions into coherent, temporally ordered sub-action units. It captures discourse-level semantics, accurately extracting motion-relevant units from abstract texts. To ensure reliable decomposition, we impose two criteria on sub-actions: (1) *Semantic completeness*, ensuring each unit forms a coherent motion; (2) *Temporal decomposability*, preserving the logical order of actions. These criteria are enforced through a structured prompt with in-context examples that guide the LLM to decompose the free-form texts. To enhance robustness, we apply LLM-based paraphrasing followed by majority voting to produce the final decomposition:

$$\mathcal{T} \rightarrow \{\mathcal{T}_1, \mathcal{T}_2, ..., \mathcal{T}_k\}$$

Further implementation details, including prompt design and voting strategies, are provided in Appendix A.2.

**Soft Masking Optimization** The second stage refines the coarse, sequence-level alignment learned during pretraining into a fine-grained, temporally localized understanding of motion. While attention pooling aggregates frame-wise features into global embeddings, it operates with a fixed receptive field over the full sequence, often suppressing local variations and blurring sub-action boundaries. This becomes particularly limiting in zero-shot settings, where unseen actions must be inferred from latent structure without supervision. To address this, SMO introduces a learnable soft mask that adapts attention pooling to each sub-action query. By dynamically reweighting frame contributions, the mask enables the model to focus on semantically relevant segments while preserving broader temporal context.

Given a set of $k$ sub-action queries $\{\mathcal{T}_1, \ldots, \mathcal{T}_k\}$, each query $\mathcal{T}_i$ is encoded as a semantic embedding $\hat{t}_i \in \mathbb{R}^d$ using the frozen $E_T$. The motion sequence $\mathcal{M} = \{m_1, \ldots, m_L\}$ is encoded by $E_M$ into frame-wise features $F = \{f_1, \ldots, f_L\}$. To identify sub-action-relevant frames, we introduce a set of learnable soft masks $[M] = \{M_1, \ldots, M_k\}$, where each $M_i \in \mathbb{R}^L$ assigns a relevance score to every frame for the $i$-th sub-action. Since each frame may be semantically relevant to multiple sub-actions, we normalize the masks across sub-actions at each frame via softmax (Liu et al. 2016):

$$\hat{M}_{i,t} = \frac{\exp(M_{i,t})}{\sum_{j=1}^{k} \exp(M_{j,t})}, \qquad (2)$$

where $\hat{M}_{i,t} \in (0, 1)$ represents the normalized probability that frame $t$ semantically belongs to sub-action $T_i$. This normalization avoids degenerate solutions (e.g., all masks attending to all frames), encourages competition between masks, and promotes segment-level disentanglement.

Each normalized soft mask $\hat{M}_i \in \mathbb{R}^L$ is used to reweigh frame-wise features via element-wise multiplication:

$$\hat{F}_i = \hat{M}_i \cdot F = \{\hat{M}_{i,1} f_1, \hat{M}_{i,2} f_2, \ldots, \hat{M}_{i,L} f_L\}. \quad (3)$$

This operation serves as a soft selection over time, allowing the model to emphasize semantically important frames while retaining the full temporal resolution.
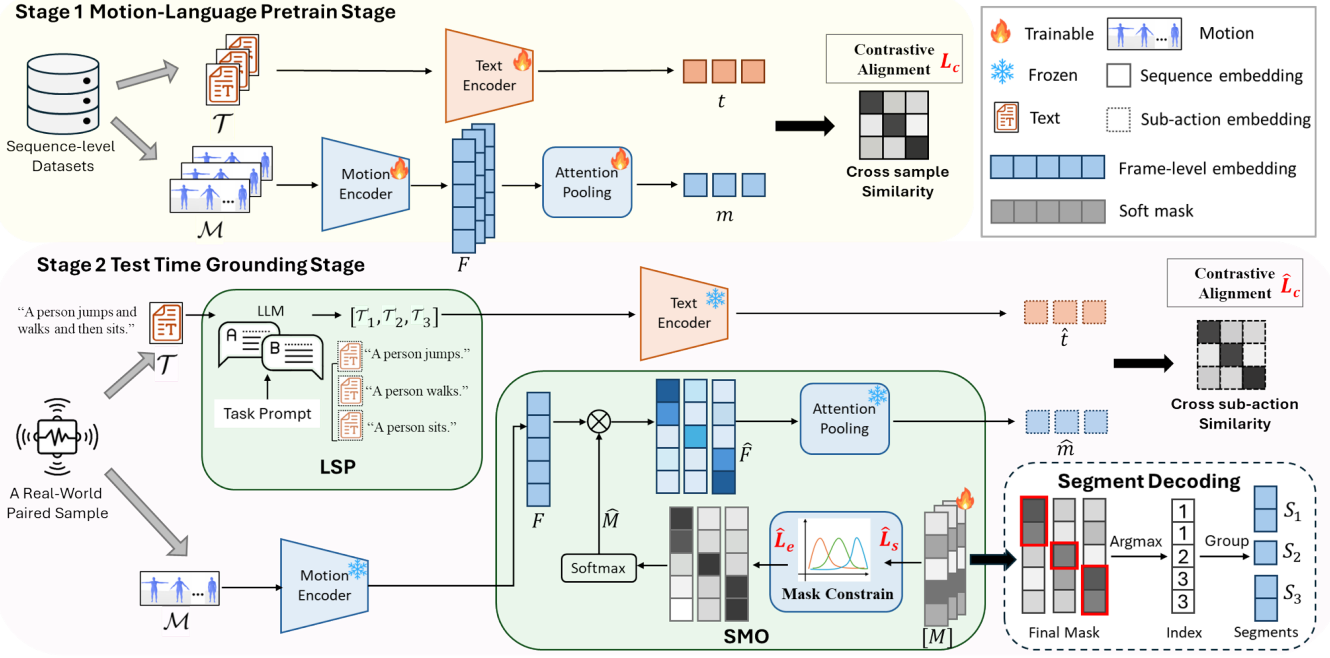
Figure 2: Overall framework of ZOMG.

Unlike hard windowing or static pooling, soft masking permits continuous gradients and supports query-specific focus without requiring discrete decisions. The resulting features $\hat{F}_i$ are then aggregated via attention pooling to produce the sub-action embedding:

$$\hat{m}_i = P(\hat{F}_i) \approx P([f_t]), \quad t \in T_i. \qquad (4)$$

Here, $\hat{M}_i$ acts as a semantic pre-filter that suppresses irrelevant frames by scaling down their feature amplitudes before pooling. This reshapes the pooling's receptive field around sub-action-relevant regions, while preserving the inter-frame dynamics learned during pretraining.

To supervise the learning of query-specific soft masks, we adopt the same contrastive loss as pretraining, encouraging each sub-action's motion embedding $\hat{m}_i$ to align with its corresponding text embedding $\hat{t}_i$. Formally:

$$\hat{\mathcal{L}}_c = -\log \frac{\exp(\hat{m}_i \cdot \hat{t}_i / \tau)}{\sum_{\hat{t} \in \{\hat{t}_i, \hat{t}^-\}} \exp(\hat{m}_i \cdot \hat{t} / \tau)}, \qquad (5)$$

where $\hat{t}^-$ are negative sub-action embeddings from the same sequence. Unlike pretraining where negatives are drawn from different samples, here we emphasize intra-sample sub-action discrimination, promoting segment-level orthogonality.

**Mask constrains.** Although soft masking enables query-specific frame reweighting, learning high-quality masks remains under-constrained. Without explicit supervision, the model may produce ambiguous or unstable masks that hinder precise localization. In particular, effective sub-action masks should satisfy two critical properties: (1) *inter-segment separability*, ensuring that different sub-actions attend to distinct temporal regions; and (2) *intra-segment continuity*, preserving the temporal coherence within each segment. These

properties are crucial for disentangling overlapping motion patterns and avoiding fragmented attention.

To promote inter-segment separability, we introduce an exclusivity loss $\hat{\mathcal{L}}_e$ that penalizes frame-level overlap between soft masks of different sub-actions, promoting temporal separation across sub-actions:

$$\hat{\mathcal{L}}_e = \frac{1}{k(k-1)} \sum_{i \neq j} \hat{M}_i^\top \hat{M}_j, \qquad (6)$$

where $\hat{M}_i^\top \hat{M}_j$ measures how much two sub-actions attend to the same frames. This constraint encourages orthogonality across masks, helping to disentangle temporally overlapping motion patterns and sharpen segment boundaries.

For intra-segment continuity, we impose a smoothness loss $\hat{\mathcal{L}}_s$ that penalizes abrupt transitions in mask values across consecutive frames:

$$\hat{\mathcal{L}}_s = \frac{1}{k(L-1)} \sum_{i=1}^{k} \sum_{t=1}^{L-1} (\hat{M}_{i,t+1} - \hat{M}_{i,t})^2. \qquad (7)$$

reflecting the temporal consistency of human motion. This regularization enforces local consistency within each mask, suppressing fragmented or unstable masks and encouraging the formation of coherent, contiguous segments.

Together, the two constraints form a complementary design that enforces both segment-level exclusivity and intra-segment continuity, enforcing meaningful structural priors over the learned soft masks. The final test-time objective integrates alignment supervision with both constraints:

$$\hat{\mathcal{L}} = \alpha \cdot \hat{\mathcal{L}}_c + \beta \cdot \hat{\mathcal{L}}_e + \gamma \cdot \hat{\mathcal{L}}_s, \qquad (8)$$
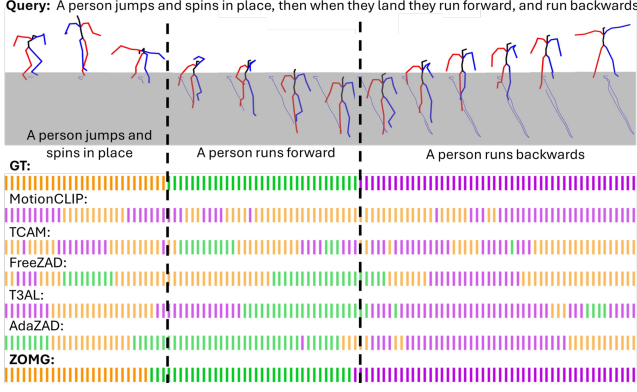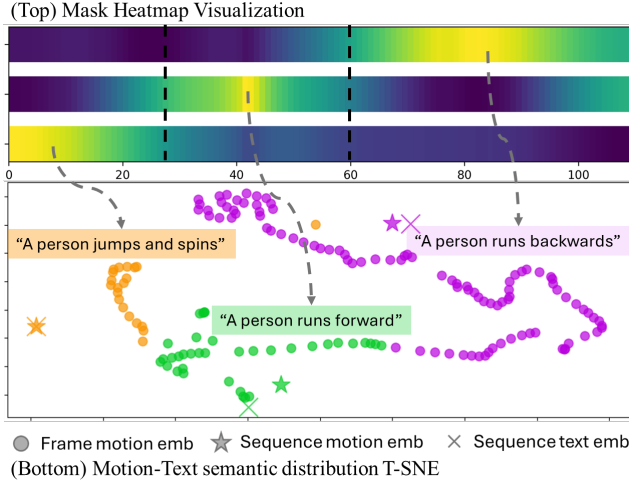
Figure 3: Motion grounding comparison in HumanML3D.



Figure 4: Analysis of ZOMG in (Top) mask heatmap, and (Bottom) T-SNE distribution.

where $\alpha, \beta, \gamma$ control the trade-off between semantic alignment, segment separability, and temporal coherence.

**Mask-Based Segment Decoding.** After obtaining the optimized soft masks, we convert them into discrete temporal segments by assigning each frame to the sub-action with the highest activation. Specifically, for each frame $t$, we compute

$$y_t = \arg\max_i \hat{M}_{i,t}, \tag{9}$$

resulting in a frame-level label sequence $y_1, \ldots, y_L$. Contiguous spans with identical labels are then grouped into segments $\{S_1, \ldots, S_k\}$. Unlike threshold-based methods that require hyperparameter tuning, our decoding is parameter-free and directly reflects model confidence, yielding semantically aligned and temporally coherent segmentation.

## Experiment

### Setup

**Dataset.** We evaluate ZOMG on three public motion-language datasets: HumanML3D (Guo et al. 2022) (29,232 motion-text pairs), KIT-ML (Plappert, Mandery, and Asfour
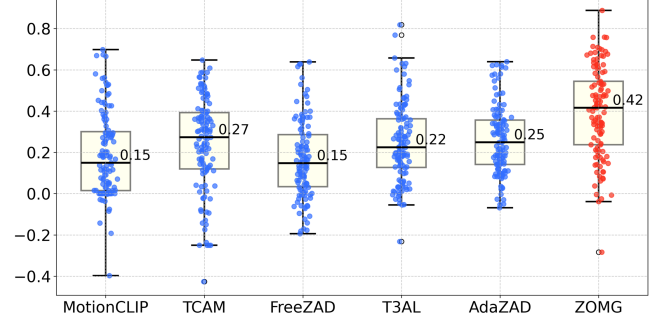


Figure 5: The boxplot comparison of semantic similarity for HumanML3D grounded motion-text pairs

2016) (3,911 pairs), and BABEL (Punnakkal et al. 2021), which includes dense annotations over 100 hours. For grounding, we select samples containing multiple sub-actions identified by the LLM, yielding 11,220 (HumanML3D), 1,207 (KIT-ML), and 2,210 (BABEL) test instances.

**Baseline.** We adopt MotionCLIP (Tevet et al. 2022) as a baseline directly computing similarity between motion and text embeddings. Due to the lack of zero-shot grounding baselines in the motion domain, we compare with SOTA video grounding methods, including non-adaptive models (TCAM (Belharbi et al. 2023), FreeZAD (Han et al. 2025)) and test-time training (TTT) methods (T3AL (Liberatori et al. 2024), AdaZAD (Han et al. 2025)). We evaluate grounding quality using mean Average Precision (mAP) computed over multiple temporal Intersection-over-Union (IoU) thresholds. Following standard practice, we report mAP averaged over thresholds of [0.3: 0.1: 0.8] on all three datasets.

**Implementation.** We adopt transformer-based encoders for motion and text. Motion is represented in H3D format with 251/263-dimensional features per frame (Lu et al. 2023). We use Qwen-Plus (Yang et al. 2025) LLM for sub-action decomposition and optimize soft masks for each query 100 steps. Loss weights for $\alpha, \beta, \gamma$ are set at 1, 0.005 and 100, respectively. All experiments are conducted on a single NVIDIA RTX 3090 GPU using PyTorch with mixed precision.

### Motion Grounding Performance

**Quantative Results.** Shown in Table 1, non-adaptive methods directly apply pretrained knowledge without instance-level adaptation, leading to sub-optimal performance due to their inability to capture fine-grained motion-language alignment. In contrast, test-time tuning approaches significantly improve grounding quality by optimizing lightweight parameters for each instance. Among them, ZOMG achieves the best performance, surpassing the previous SOTA method AdaZAD by +8.69%, +9.62%, and +3.29% mAP on HumanML3D, KIT-ML, and BABEL, respectively. These consistent gains validate the robustness and effectiveness of ZOMG in producing fine-grained, semantically aligned motion segments under zero-shot open-vocabulary conditions.

**Visualization Results.** Figure 3 presents qualitative comparisons across different methods. While baselines often pro-

| Dataset | Method | TTT | AP@8↑ | AP@7↑ | AP@6↑ | AP@5↑ | AP@4↑ | AP@3↑ | mAP ↑ |
|---|---|---|---|---|---|---|---|---|---|
| HumanML3D | MotionCLIP | ✗ | 0.00 | 2.13 | 4.61 | 12.06 | 17.85 | 24.29 | 10.16 |
| | TCAM | ✗ | 1.06 | 5.14 | 13.12 | 19.92 | 31.09 | 36.58 | 17.82 |
| | FreeZAD | ✗ | 7.68 | 14.17 | 21.09 | 27.24 | 34.81 | 42.76 | 24.63 |
| | T3AL | ✓ | 13.71 | 22.64 | 34.57 | 46.99 | 55.14 | 67.08 | 40.21 |
| | AdaZAD | ✓ | 15.11 | 24.39 | 38.42 | 48.18 | 56.87 | 67.63 | 41.77 |
| | ZOMG | ✓ | **28.25** | **37.53** | **51.42** | **53.90** | **60.85** | **70.83** | **50.46** |
| KIT-ML | MotionCLIP | ✗ | 0.00 | 1.58 | 3.72 | 9.62 | 17.31 | 23.93 | 9.36 |
| | TCAM | ✗ | 3.85 | 7.54 | 12.85 | 21.10 | 28.46 | 35.23 | 18.17 |
| | FreeZAD | ✗ | 6.79 | 9.34 | 14.39 | 22.62 | 27.43 | 37.50 | 19.68 |
| | T3AL | ✓ | 11.54 | 15.38 | 18.29 | 27.78 | 36.97 | 54.91 | 27.48 |
| | AdaZAD | ✓ | 13.68 | 17.91 | 24.57 | 32.14 | 39.49 | 56.88 | 30.78 |
| | ZOMG | ✓ | **22.01** | **25.64** | **39.69** | **43.74** | **48.59** | **62.73** | **40.40** |
| BABEL | MotionCLIP | ✗ | 0.00 | 1.00 | 1.49 | 3.98 | 7.96 | 12.44 | 4.31 |
| | TCAM | ✗ | 0.00 | 1.00 | 2.49 | 4.98 | 11.44 | 17.91 | 6.30 |
| | FreeZAD | ✗ | 0.00 | 1.32 | 2.56 | 5.97 | 10.95 | 17.31 | 6.35 |
| | T3AL | ✓ | 0.23 | 1.86 | 3.41 | 8.09 | 16.38 | 29.57 | 9.92 |
| | AdaZAD | ✓ | **0.72** | **2.64** | 4.93 | 11.82 | 18.70 | 32.09 | 10.15 |
| | ZOMG | ✓ | 0.50 | 2.49 | **5.26** | **13.16** | **23.03** | **36.18** | **13.44** |

Table 1: Motion grounding results on three public datasets. (TTT denotes Test-Time Training.)

| TTT | LSP | SMO | | | AP@7 | AP@3 | mAP |
|---|---|---|---|---|---|---|---|
| | | [M] | $\hat{\mathcal{L}}_s$ | $\hat{\mathcal{L}}_e$ | | | |
| ✗ | ✗ | ✗ | ✗ | ✗ | 0.91 | 14.33 | 5.86 |
| ✗ | ✓ | ✗ | ✗ | ✗ | 2.13 | 24.29 | 10.16 |
| ✓ | ✓ | ✗ | ✗ | ✗ | 13.21 | 38.52 | 22.69 |
| ✓ | ✓ | ✓ | ✗ | ✗ | 21.17 | 62.90 | 38.63 |
| ✓ | ✓ | ✓ | ✓ | ✗ | 25.76 | 66.29 | 41.54 |
| ✓ | ✓ | ✓ | ✗ | ✓ | 33.51 | **71.20** | 47.28 |
| ✓ | ✓ | ✓ | ✓ | ✓ | **37.53** | 70.83 | **50.46** |

Table 2: Ablation study of ZOMG on HumanML3D.

duce fragmented or misaligned segments, ZOMG generates temporally coherent and semantically meaningful outputs that align well with the intended sub-actions. Notably, the predicted segments of ZOMG closely match the optimized mask heatmaps in Figure 4 (top), confirming that soft masking effectively captures subtle temporal structure and aligns motion with nuanced open-vocabulary semantics. Additional visualizations are provided in the Appendix A.4.

**Semantic Alignment Evaluation.** To further assess grounding quality, we evaluate the semantic alignment between the segmented motion clips and their corresponding texts. After applying different methods to segment the motion into sub-actions, we compute the motion-text similarity for each pair using a pretrained motion-language encoder. This metric serves as an intuitive proxy for alignment quality: higher similarity indicates that the segmented motion more faithfully reflects the intended semantics. As shown in Figure 5, ZOMG consistently achieves stronger motion-text alignment, reflected both in its top-ranking similarity scores and in a distribution skewed toward higher values. This alignment-based evaluation highlights ZOMG's ability to produce semantically faithful segments, offering clear

advantages in real-world, open-vocabulary scenarios.

## Ablation and Analytical Studies

**Component Ablation.** Table 2 assess the contribution of each component in ZOMG's test-time grounding stage on HumanML3D. Without test-time training, the non-adaptive baseline performs poorly (5.86% mAP), highlighting the need for instance-specific optimization. Incorporating LSP yields clear gains over rule-based segmentation, demonstrating the value of LLM-guided sub-action decomposition. To evaluate SMO, we progressively activate its components. Introducing the soft mask [M] significantly boosts performance, suggesting that optimizing frame-wise weights is critical for aligning motion with language. Adding $\hat{\mathcal{L}}_s$ or $\hat{\mathcal{L}}_e$ further improves results, with the full model achieving the best performance (50.46% mAP). These findings confirm the complementary roles of the mask constraints and the importance of each module for effective zero-shot grounding.

| Method | Param.↓ | GFLOPs↓ | Samples/s↑ | mAP↑ |
|---|---|---|---|---|
| T3AL | 1.2 M | 2528.1 | 6.87 | 40.21 |
| AdaZAD | 1.2 M | 2675.9 | 6.65 | 41.77 |
| ZOMG | **0.5 K** | **302.5** | **23.25** | **50.46** |

Table 3: Computational costs during TTT for 100 steps.

**Mask Quality Analysis.** To evaluate the learned masks, we examine two key properties: inter-segment separation and intra-segment continuity, both essential for discovering temporally distinct and semantically coherent sub-actions. **Qualitatively**, Figure 4 (top) shows that the optimized masks activate over localized regions with smooth transitions and minimal overlap. Each mask captures a distinct sub-action with clear and interpretable boundaries, dynamically adjusting the receptive field of pooling layers to focus on the most relevant frames. **Quantitatively**, We further assess mask quality using

| Dataset | Protocol | Task | SOTA methods | | | | Augmentation methods | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | | MotionCLIP | TEMOS | TMR | MESM | Noise | Scaling | Concat | ZOMG |
| **HumanML3D** | A | T2M | 16.00 | 13.15 | 16.32 | 19.20 | 17.05 | 17.15 | 17.15 | **19.38** |
| | | M2T | 16.95 | 7.74 | 18.72 | **21.30** | 17.98 | 18.57 | 18.58 | 20.53 |
| | B | T2M | 22.49 | 12.36 | 22.75 | 26.29 | 23.87 | 23.96 | 23.90 | **27.01** |
| | | M2T | 21.69 | 9.96 | 23.25 | 25.09 | 23.19 | 23.61 | 23.71 | **25.52** |
| | C | T2M | 86.90 | 62.05 | 84.42 | 86.22 | 85.92 | 86.27 | 85.89 | **87.58** |
| | | M2T | 87.12 | 62.71 | 84.50 | 86.15 | 86.18 | 86.72 | 85.93 | **87.56** |
| **KIT-ML** | A | T2M | 22.49 | 19.54 | 22.00 | 23.76 | 23.19 | 23.61 | 23.86 | **24.03** |
| | | M2T | 22.92 | 22.03 | 22.36 | 21.30 | 23.65 | 22.54 | 23.15 | **23.71** |
| | B | T2M | 40.22 | 34.48 | 41.52 | **43.04** | 39.66 | 39.15 | 40.07 | 42.04 |
| | | M2T | 34.10 | 30.29 | 34.02 | 35.61 | 34.69 | 33.57 | 36.80 | **37.42** |
| | C | T2M | 75.68 | 65.58 | 76.03 | 76.11 | 76.28 | 75.53 | 76.70 | **77.43** |
| | | M2T | 75.61 | 64.88 | 75.55 | 76.65 | 75.94 | 75.80 | 75.80 | **76.99** |

Table 4: Text-to-motion retrieval performance comparison.

grounding error and continuity rate over TTT iterations (Figure 6). In (a), ZOMG consistently reduces grounding error across iterations and outperforms all baselines. Notably, this performance improvement closely follows the decline of the exclusivity loss $\hat{\mathcal{L}}_e$, indicating that enhanced inter-mask separation directly benefits grounding accuracy. In (b), continuity rate quickly declines and remains low, reflecting the effectiveness of the smoothness constraint $\hat{\mathcal{L}}_s$. Together, these results confirm that ZOMG learns structured and disentangled masks through efficient test-time adaptation.

**Semantic Prior Validation.** ZOMG assumes that pretrained frame-level representations already encode meaningful temporal semantics, which can be further refined through instance-wise optimization. To validate this prior, we visualize the embeddings using t-SNE. As shown in Figure 4 (bottom), frame embeddings form smooth trajectories within each sub-action, indicating coherent transitions and reflecting temporal semantics. This confirms that the pretrained space provides a strong structural prior, which our instance-specific optimization can exploit for fine-grained motion grounding.

**Efficiency and Temporal Scaling.** Table 3 shows ZOMG delivers strong grounding performance with minimal test-time cost. By optimizing only soft masks while keeping the encoder frozen, it achieves over 3× higher throughput than existing TTT methods, enabling practical deployment in latency-sensitive, annotation-free scenarios. Fig. 6 (a) further exhibits favorable *temporal scaling*, where error steadily decreases with more optimization steps. This controllable trade-off between adaptation time and accuracy allows flexible use under different computational budgets.

### Downstream Benefits

To assess the broader utility of ZOMG, we augment motion retrieval datasets by generating fine-grained motion-text pairs grounded from sub-action queries, enriching compositional semantics beyond original annotations. We evaluate
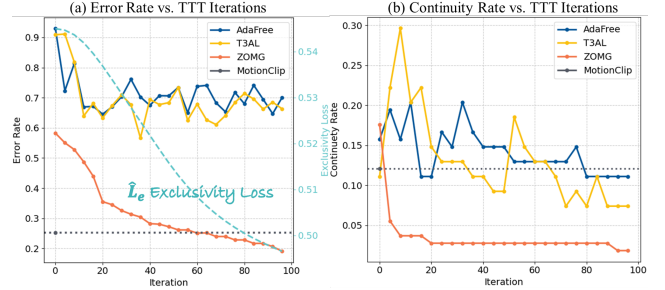


Figure 6: Comparison of (a) Mask constrain performance and (b) Motion-text semantic similarity on HumanML3D.

their impact on SOTA methods (MotionCLIP, TEMOS, TMR, MESM (Shi and Zhang 2024)) under three standard protocols, including filtered and batch retrieval settings. For fair comparison, we include several augmentation baselines such as noise injection, temporal scaling, and motion concatenation, with all methods producing an equal number of samples. As shown in Table 4, ZOMG consistently improves retrieval across models and settings, indicating that our grounded segments are both temporally precise and semantically discriminative. Full implementation details are in Appendix A.5.

## Conclusion

Our ZOMG demonstrates that accurate and open-vocabulary motion grounding can be achieved annotation-free, through lightweight test-time optimization without modifying pretrained models. Its high efficiency and strong performance across grounding and retrieval tasks highlight the potential of instance-adaptive inference for real-world deployment. By uncovering compositional motion units in an unsupervised manner, ZOMG provides a scalable foundation for interpretable motion understanding and broad downstream transfer.

## Acknowledgments

## References

Athanasiou, N.; Petrovich, M.; Black, M. J.; and Varol, G. 2022. Teach: Temporal action composition for 3d humans. In *2022 International Conference on 3D Vision (3DV)*, 414–423. IEEE.

Belharbi, S.; Ben Ayed, I.; McCaffrey, L.; and Granger, E. 2023. Tcam: Temporal class activation maps for object localization in weakly-labeled unconstrained videos. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 137–146.

Bordes, F.; Pang, R. Y.; Ajay, A.; Li, A. C.; Bardes, A.; Petryk, S.; Mañas, O.; Lin, Z.; Mahmoud, A.; Jayaraman, B.; et al. 2024. An introduction to vision-language modeling. *arXiv preprint arXiv:2405.17247*.

Chao, Y.-W.; Vijayanarasimhan, S.; Seybold, B.; Ross, D. A.; Deng, J.; and Sukthankar, R. 2018. Rethinking the faster r-cnn architecture for temporal action localization. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 1130–1139.

Chen, Y.-W.; Tsai, Y.-H.; and Yang, M.-H. 2021. End-to-end multi-modal video temporal grounding. *Advances in Neural Information Processing Systems*, 34: 28442–28453.

FitzGerald, N.; Michael, J.; He, L.; and Zettlemoyer, L. 2018. Large-scale QA-SRL parsing. *arXiv preprint arXiv:1805.05377*.

Gao, J.; Sun, C.; Yang, Z.; and Nevatia, R. 2017. Tall: Temporal activity localization via language query. In *Proceedings of the IEEE international conference on computer vision*, 5267–5275.

Gorti, S. K.; Vouitsis, N.; Ma, J.; Golestan, K.; Volkovs, M.; Garg, A.; and Yu, G. 2022. X-pool: Cross-modal language-video attention for text-video retrieval. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 5006–5015.

Guo, C.; Zou, S.; Zuo, X.; Wang, S.; Ji, W.; Li, X.; and Cheng, L. 2022. Generating diverse and natural 3d human motions from text. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 5152–5161.

Han, C.; Wang, H.; Kuang, J.; Zhang, L.; and Gui, J. 2025. Training-Free Zero-Shot Temporal Action Detection with Vision-Language Models. *arXiv preprint arXiv:2501.13795*.

Kong, Q.; Wu, Z.; Deng, Z.; Klinkigt, M.; Tong, B.; and Murakami, T. 2019. Mmact: A large-scale dataset for cross modal human action understanding. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 8658–8667.

Lee, J.; Lee, I.; and Kang, J. 2019. Self-attention graph pooling. In *International conference on machine learning*, 3734–3743. pmlr.

Li, K.; and Feng, Y. 2024. Motion generation from fine-grained textual descriptions. *arXiv preprint arXiv:2403.13518*.

Liberatori, B.; Conti, A.; Rota, P.; Wang, Y.; and Ricci, E. 2024. Test-time zero-shot temporal action localization. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 18720–18729.

Lin, J.; Zeng, A.; Lu, S.; Cai, Y.; Zhang, R.; Wang, H.; and Zhang, L. 2023. Motion-x: A large-scale 3d expressive whole-body human motion dataset. *Advances in Neural Information Processing Systems*, 36: 25268–25280.

Liu, D.; Qu, X.; Zhou, P.; and Liu, Y. 2022. Exploring motion and appearance information for temporal sentence grounding. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, 1674–1682.

Liu, W.; Wen, Y.; Yu, Z.; and Yang, M. 2016. Large-margin softmax loss for convolutional neural networks. *arXiv preprint arXiv:1612.02295*.

Liu, Y.; Kothari, P.; Van Delft, B.; Bellot-Gurlet, B.; Mordan, T.; and Alahi, A. 2021. Ttt++: When does self-supervised test-time training fail or thrive? *Advances in Neural Information Processing Systems*, 34: 21808–21820.

Lu, S.; Chen, L.-H.; Zeng, A.; Lin, J.; Zhang, R.; Zhang, L.; and Shum, H.-Y. 2023. Humantomato: Text-aligned whole-body motion generation. *arXiv preprint arXiv:2310.12978*.

Nag, S.; Zhu, X.; Song, Y.-Z.; and Xiang, T. 2022. Zero-shot temporal action detection via vision-language prompting. In *European conference on computer vision*, 681–697. Springer.

Nguyen, T. T.; Bin, Y.; Wu, X.; Hu, Z.; Nguyen, C.-D. T.; Ng, S.-K.; and Luu, A. T. 2025. Multi-scale contrastive learning for video temporal grounding. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 39, 6227–6235.

Petrovich, M.; Black, M. J.; and Varol, G. 2022. Temos: Generating diverse human motions from textual descriptions. In *European Conference on Computer Vision*, 480–497. Springer.

Petrovich, M.; Black, M. J.; and Varol, G. 2023. Tmr: Text-to-motion retrieval using contrastive 3d human motion synthesis. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 9488–9497.

Plappert, M.; Mandery, C.; and Asfour, T. 2016. The kit motion-language dataset. *Big data*, 4(4): 236–252.

Punnakkal, A. R.; Chandrasekaran, A.; Athanasiou, N.; Quiros-Ramirez, A.; and Black, M. J. 2021. BABEL: Bodies, action and behavior with english labels. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 722–731.

Qian, J.; Zhu, Z.; Zhou, H.; Feng, Z.; Zhai, Z.; and Mao, K. 2025. Beyond the Next Token: Towards Prompt-Robust Zero-Shot Classification via Efficient Multi-Token Prediction. *arXiv preprint arXiv:2504.03159*.

Radford, A.; Kim, J. W.; Hallacy, C.; Ramesh, A.; Goh, G.; Agarwal, S.; Sastry, G.; Askell, A.; Mishkin, P.; Clark, J.; et al. 2021. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, 8748–8763. PmLR.

Shi, H.; and Zhang, H. 2024. Modal-Enhanced Semantic Modeling for Fine-Grained 3D Human Motion Retrieval. In *Proceedings of the 32nd ACM International Conference on Multimedia*, 10114–10123.

Sun, Y.; Wang, X.; Liu, Z.; Miller, J.; Efros, A.; and Hardt, M. 2020. Test-time training with self-supervision for generalization under distribution shifts. In *International conference on machine learning*, 9229–9248. PMLR.

Tevet, G.; Gordon, B.; Hertz, A.; Bermano, A. H.; and Cohen-Or, D. 2022. Motionclip: Exposing human motion generation to clip space. In *European Conference on Computer Vision*, 358–374. Springer.

Wang, J.; Lin, X.; Huang, H.; Ke, X.; Wu, R.; You, C.; and Guo, K. 2023a. GLANet: temporal knowledge graph completion based on global and local information-aware network. *Applied Intelligence*, 53(16): 19285–19301.

Wang, M.; Xing, J.; Mei, J.; Liu, Y.; and Jiang, Y. 2023b. Actionclip: Adapting language-image pretrained models for video action recognition. *IEEE Transactions on Neural Networks and Learning Systems*.

Wang, R.; Tang, D.; Duan, N.; Wei, Z.; Huang, X.; Cao, G.; Jiang, D.; Zhou, M.; et al. 2020. K-adapter: Infusing knowledge into pre-trained models with adapters. *arXiv preprint arXiv:2002.01808*.

Wang, X.; Kang, Z.; and Mu, Y. 2024. Text-controlled motion mamba: text-instructed temporal grounding of human motion. *arXiv preprint arXiv:2404.11375*.

Wortsman, M.; Ilharco, G.; Kim, J. W.; Li, M.; Kornblith, S.; Roelofs, R.; Lopes, R. G.; Hajishirzi, H.; Farhadi, A.; Namkoong, H.; et al. 2022. Robust fine-tuning of zero-shot models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 7959–7971.

Wu, L.-T.; Lin, J.-R.; Leng, S.; Li, J.-L.; and Hu, Z.-Z. 2022. Rule-based information extraction for mechanical-electrical-plumbing-specific semantic web. *Automation in Construction*, 135: 104108.

Xu, H.; Ghosh, G.; Huang, P.-Y.; Arora, P.; Aminzadeh, M.; Feichtenhofer, C.; Metze, F.; and Zettlemoyer, L. 2021. Vlm: Task-agnostic video-language model pre-training for video understanding. *arXiv preprint arXiv:2105.09996*.

Xue, W.; Qian, C.; Wu, J.; Zhou, Y.; Liu, W.; Ren, J.; Fan, S.; and Zhang, Y. 2025. ShotVL: Human-Centric Highlight Frame Retrieval via Language Queries. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 39, 9050–9058.

Yan, S.; Xiong, X.; Nagrani, A.; Arnab, A.; Wang, Z.; Ge, W.; Ross, D.; and Schmid, C. 2023. Unloc: A unified framework for video localization tasks. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 13623–13633.

Yang, A.; Li, A.; Yang, B.; Zhang, B.; Hui, B.; Zheng, B.; Yu, B.; Gao, C.; Huang, C.; Lv, C.; et al. 2025. Qwen3 Technical Report. *arXiv preprint arXiv:2505.09388*.

Yang, A.; Miech, A.; Sivic, J.; Laptev, I.; and Schmid, C. 2022a. Tubedetr: Spatio-temporal video grounding with transformers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 16442–16453.

Yang, J.; Zhou, Y.; Huang, H.; Zou, H.; and Xie, L. 2022b. MetaFi: Device-free pose estimation via commodity WiFi for metaverse avatar simulation. In *2022 IEEE 8th World Forum on Internet of Things (WF-IoT)*, 1–6. IEEE.

Yu, S.; Cho, J.; Yadav, P.; and Bansal, M. 2023. Self-chained image-language model for video localization and question answering. *Advances in Neural Information Processing Systems*, 36: 76749–76771.

Zeng, L.-A.; Huang, G.; Wu, G.; and Zheng, W.-S. 2025. Light-t2m: A lightweight and fast model for text-to-motion generation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 39, 9797–9805.

Zhang, M.; Li, H.; Cai, Z.; Ren, J.; Yang, L.; and Liu, Z. 2023. Finemogen: Fine-grained spatio-temporal motion generation and editing. *Advances in Neural Information Processing Systems*, 36: 13981–13992.

Zhang, S.; Peng, H.; Fu, J.; and Luo, J. 2020. Learning 2d temporal adjacent networks for moment localization with natural language. In *Proceedings of the AAAI conference on artificial intelligence*, volume 34, 12870–12877.

Zhao, W. X.; Zhou, K.; Li, J.; Tang, T.; Wang, X.; Hou, Y.; Min, Y.; Zhang, B.; Zhang, J.; Dong, Z.; et al. 2023. A survey of large language models. *arXiv preprint arXiv:2303.18223*, 1(2).

Zhao, Y.; Xiong, Y.; Wang, L.; Wu, Z.; Tang, X.; and Lin, D. 2017. Temporal action detection with structured segment networks. In *Proceedings of the IEEE international conference on computer vision*, 2914–2923.

Zhou, Y.; Huang, H.; Yuan, S.; Zou, H.; Xie, L.; and Yang, J. 2023a. MetaFi++: WiFi-enabled transformer-based human pose estimation for metaverse avatar simulation. *IEEE Internet of Things Journal*, 10(16): 14128–14136.

Zhou, Y.; Yang, J.; Huang, H.; and Xie, L. 2024. Adapose: Towards cross-site device-free human pose estimation with commodity wifi. *IEEE Internet of Things Journal*.

Zhou, Y.; Yang, J.; Zou, H.; and Xie, L. 2023b. Tent: Connect language models with iot sensors for zero-shot activity recognition. *arXiv preprint arXiv:2311.08245*.

Zhou, Z.; Wan, Y.; and Wang, B. 2024. Avatargpt: All-in-one framework for motion understanding planning generation and beyond. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 1357–1366.