# RECRUITVIEW: A Multimodal Dataset for Predicting Personality and Interview Performance for Human Resources Applications

Amit Kumar Gupta[1][*][†]    Farhan Sheth[1][*]    Hammad Shaikh[1]    Dheeraj Kumar[1]

Angkul Puniya[1]    Deepak Panwar[1]    Sandeep Chaurasia[1]    Priya Mathur[2]

[1]Manipal University Jaipur, India    [2]Poornima Institute of Engineering & Technology, India

{amit.gupta, deepak.panwar, sandeep.chaurasia}@jaipur.manipal.edu
{farhan.219310185, hammad.229301534}@muj.manipal.edu
{dheeraj.229301593, angkul.23FE10CAI00309}@muj.manipal.edu
priya.mathur@poornima.org

## Abstract

*Automated personality and soft skill assessment from multimodal behavioral data remains challenging due to limited datasets and methods that fail to capture geometric structure inherent in human traits. We introduce RE-CRUITVIEW, a dataset of 2,011 naturalistic video interview clips from 300+ participants with 27,000 pairwise comparative judgments across 12 dimensions: Big Five personality traits, overall personality score, and six interview performance metrics. To leverage this data, we propose Cross-Modal Regression with Manifold Fusion (CRMF), a geometric deep learning framework that explicitly models behavioral representations across hyperbolic, spherical, and Euclidean manifolds. CRMF employs geometry-specific expert networks to capture hierarchical trait structures, directional behavioral patterns, and continuous performance variations simultaneously. An adaptive routing mechanism dynamically weights expert contributions based on input characteristics. Through principled tangent space fusion, CRMF achieves superior performance while training 40–50% fewer trainable parameters than large multimodal models. Extensive experiments demonstrate that CRMF substantially outperforms the selected baselines, achieving up to 11.4% improvement in Spearman correlation and 6.0% in concordance index. Our RECRUITVIEW dataset is publicly available at https://huggingface.co/datasets/AI4A-lab/RecruitView.*

## 1. Introduction

Interviews are integral to hiring, coaching, and clinical evaluation. Judgments hinge on subtle behaviors distributed across what candidates say, how they speak, and how they present visually. As video interviewing scales, computational assessment must read these signals coherently rather than in isolation. Estimating personality traits and interview performance from short video responses is a multimodal problem spanning vision, speech, and language. Interviews rely on complementary lexical, prosodic, and visual cues, therefore computational models must capture these complementary signals without discarding their structure.

General-purpose LMMs (MiniCPM-o [1], VideoL-LaMA2 [2], Qwen2.5-Omni [3]) offer breadth but are not tuned for fine-grained social inference. Conventional fusion maps all modalities to a single Euclidean latent via concatenation or vanilla attention, ignoring modality-specific geometry. A single latent geometry limits representational adequacy.

Progress is also constrained by supervision: existing datasets are noisy, weakly controlled, and often not domain-specific; they typically lack multi-trait personality and interview-related metrics, instead relying on direct scalar ratings that are sensitive to scale use and inter-rater variability. To address this, we introduce RECRUITVIEW—**R**ecorded **E**valuations of **C**andidate **R**esponses for **U**nderstanding **I**ndividual **T**raits—a multimodal interview corpus of 2,011 clips from more than 300 sessions, each aligned to one of 76 questions. Clinical psychologists provided about 27,000 pairwise comparisons between answers to the same prompt, which we convert into continuous scores for 12 targets, namely the Big Five traits [4], an overall personality score, and six interview performance metrics, using a nuclear-norm-regularized multinomial logit model. This protocol reduces rater calibration

biases and yields reliable regression labels.

We propose **C**ross-Modal **R**egression with **M**anifold **F**usion (CRMF), a geometry-aware framework that projects fused multimodal features to hyperbolic, spherical, and Euclidean spaces, processes each with a geometry-specific expert, and aggregates them through input-adaptive routing with geometry-aware attention and tangent-space fusion. This design preserves manifold consistency while enabling input-conditioned combination for multi-target regression.

**The contribution of this work** is fourfold: (i) RE-CRUITVIEW, a multimodal interview dataset with psychometrically grounded labels derived from pairwise judgments mapped to continuous scores, covering 12 targets across personality and performance; (ii) CRMF, a principled geometry-aware fusion framework that learns in hyperbolic, spherical, and Euclidean spaces; (iii) an adaptive routing and geometry-aware attention mechanism with tangent-space fusion for input-conditioned combination of geometric experts; and (iv) a comprehensive evaluation demonstrating consistent gains over recent LMM baselines on all metrics.

## 2. Related Works

### 2.1. Personality and Behavioral Assessment

Automated personality and performance assessment has relied on datasets such as ChaLearn [5] and POM [6], which advanced the field but remain limited by controlled settings and narrow labeling scopes. Later efforts like YouTube Personality [7] and Interview2Personality [8] moved toward more naturalistic or interview-style data yet still suffer from smaller scale, scripted responses, and subjective absolute ratings. However, RECRUITVIEW is an in-the-wild interview dataset, where labels are obtained via pairwise comparisons, yielding consistent continuous scores. Beyond the Big Five and an overall personality index, RECRUITVIEW annotates performance dimensions (e.g., confidence, communication), enabling joint modeling of personality and interview behavior.

### 2.2. Multimodal Fusion for Behavioral Analysis

Multimodal fusion has been extensively studied for affective computing and personality recognition. Early work focused on feature-level concatenation [9] or attention-based aggregation [10, 11]. Transformer-based architectures have recently dominated this space, with methods like MULT [10] employing cross-modal attention for temporal alignment. However, these approaches operate entirely in Euclidean space, potentially missing important geometric structure in behavioral data.

Recent large multimodal models have shown remarkable zero-shot and few-shot capabilities. MiniCPM-o [1] employs an end-to-end training paradigm with modality-adaptive modules, while VideoLLaMA2 [2] introduces spatial-temporal visual token compression for efficient video understanding. Qwen2.5-Omni [3] extends text-centric LLMs with native audio-visual understanding through cross-attention fusion. Despite their general-purpose success, these models lack task-specific inductive biases for personality assessment and are not optimized for capturing the geometric properties of behavioral traits.

### 2.3. Geometric Deep Learning

Geometric deep learning extends neural networks to non-Euclidean domains. Hyperbolic neural networks [12, 13] leverage the exponentially growing capacity of hyperbolic space to model hierarchical data, showing benefits for tree-structured tasks and knowledge graph reasoning. Spherical networks [14, 15] operate on the unit sphere, naturally suited for directional data and rotational equivariance. Recent work has explored mixed-curvature spaces [16, 17] that combine multiple geometries, though primarily for representation learning rather than multimodal fusion.

Manifold-valued neural networks [18, 19] perform operations directly on Riemannian manifolds, ensuring geometric consistency. However, these methods have seen limited application in behavioral analysis. Our work is the first to systematically leverage multiple geometric manifolds for multimodal behavioral assessment with learned adaptive fusion.

### 2.4. Mixture-of-Experts Architectures

Mixture-of-experts (MoE) models [20, 21] decompose complex tasks into specialized sub-networks selected by a gating function. Traditional MoE aims for sparse activation to increase model capacity efficiently. Recent work has extended MoE to multimodal settings [22, 23] and to geometric spaces [19]. However, existing geometric MoE methods typically focus on sparsity for computational efficiency rather than complementary geometric reasoning. Our routing mechanism differs fundamentally: rather than encouraging specialization, we promote diversity to leverage complementary geometric views of behavioral data, with all experts contributing to the final prediction through learned weighting.

## 3. RECRUITVIEW

To satisfy the critical necessity of a psychometrically robust dataset to analyze multimodal interview performance and personality, we introduce RECRUITVIEW. This novel dataset comprises 2,011 video segments sampled from 331 distinct interview sessions. Specifically developed to facilitate training and testing on demanding, human-centered traits, it provides robust, continuous labels on 12 distinct targets. The dataset's key contribution is in the form of an

annotation method through pairwise ratings by clinical psychologists to mitigate rater bias and provide robust, continuous scores. The following sections describe the deliberate process of developing it in stages from stimulus creation to final form of data.

## 3.1. Data Collection

The creation of RECRUITVIEW followed a two-phase approach: creation of a broad-based question repository to be used as a prompter, and the procurement of video replies by a diverse group of respondents.

Our dataset has as its base a specially selected pool of queries crafted to elicit responses suitable to human resources evaluation and personality rating. For this purpose, we carried out an exhaustive compilation exercise from a variety of sources. We went through public domain material on interview preparation by market leaders, analyzed frequently posed queries on professional networking platforms, and closely interacted with clinical psychologists. This made the queries relevant not only in the context of professional hiring but also well suited to exploring the underlying Big Five personality dimensions.

This procedure gave rise to 76 standard interview questions (the full list of questions is available in Appendix A.1). The questions were then systematically and in a balanced manner sorted into 15 individual sets for convenience in the collection of data. There were five individual questions in each set and a typical opening question ("Introduce yourself" or "Tell me about yourself"), thereby establishing a comparable baseline in the majority of interviews but facilitating greater query variety within the dataset.

Participants were students from various Manipal universities, who responded via a custom web platform and were randomly assigned to one of 15 question sets. Interviews were recorded in diverse in-the-wild settings (e.g., classrooms, private residences). Implementation details of the platform are provided in Appendix A.2.

Data collection and use followed institutional ethical approval processes; detailed ethics, consent, and risk-mitigation discussion is provided in Appendix B.

## 3.2. Dataset Annotation

### 3.2.1. Annotation Protocol

To ensure consistency and psychometric reliability in subjective evaluations, we employed a pairwise comparison protocol inspired by prior multimodal labeling frameworks such as the ChaLearn dataset [5, 24]. Instead of assigning absolute scores, clinical psychologists were presented with two clips responding to the same interview question and asked to identify which participant better demonstrated a target attribute, for example, "Who appears more confident?" Annotators could also indicate a tie when both clips were judged equivalent. This comparative design mini-mizes calibration bias, reduces inter-rater variability, and enhances reliability in perceptual assessments [25]. The protocol was applied across twelve target dimensions covering the Big Five personality traits (Openness, Conscientiousness, Extraversion, Agreeableness, and Neuroticism), Overall Personality, and six interview or performance-related metrics: Interview Score, Answer Score, Speaking Skills, Confidence Score, Facial Expression, and Overall Performance. In total, approximately 27,310 pairwise judgments were collected, forming the basis for deriving continuous and psychometrically grounded labels.

### 3.2.2. Model Selection and Multinomial Logit (MNL)

We evaluated several frameworks for converting pairwise judgments into continuous labels, including Elo rating [26], Bradley-Terry-Luce (BTL) [27, 28], TrueSkill [29], and Glicko-2 [30]. While these models are widely used in ranking applications, they either assume strong independence across traits or lack convex formulations with clear identifiability guarantees. After empirical comparison (results in Appendix A.3), we selected the Multinomial Logit (MNL) [31] model with nuclear norm regularization, which offered both strong theoretical grounding and robust empirical performance on our dataset.

Each video $j$ is associated with a latent utility $\theta_j \in \mathbb{R}$. For a comparison between clips $j_1$ and $j_2$, the MNL model defines the probability that $j_1$ is preferred as

$$\Pr\{j_1 \succ j_2\} = \frac{\exp(\theta_{j_1})}{\exp(\theta_{j_1}) + \exp(\theta_{j_2})} \quad (1)$$

Letting $X^{(i)}$ denote the design matrix for comparison $i$ and $y_i \in \{0, 1\}$ its observed outcome, the normalized log-likelihood across $n$ comparisons is

$$\mathcal{L}(\Theta) = \frac{1}{n} \sum_{i=1}^{n} \left( y_i \langle \Theta, X^{(i)} \rangle - \log \left( 1 + \exp(\langle \Theta, X^{(i)} \rangle) \right) \right)$$
$$(2)$$

where $\Theta \in \mathbb{R}^{N \times T}$ is the matrix of utilities across $N$ videos and $T = 12$ targets.

### 3.2.3. Nuclear Norm Regularization and Optimization

Recovering utilities requires regularization to address limited sampling and correlations across traits. We therefore estimate $\Theta$ by solving the convex program

$$\hat{\Theta} = \arg \min_{\Theta \in \Omega} \left[ -\alpha \, \mathcal{L}(\Theta) + \lambda \| L^{1/2} \Theta \|_* \right] \quad (3)$$

where $\| \cdot \|_*$ denotes the nuclear norm [32], $L$ is the Laplacian of the comparison graph, and $\Omega$ constrains identifiability (e.g., centering utilities). The Laplacian-induced nuclear norm encourages low-rank structure while respecting the blockwise nature of the pairwise comparisons (same-question groups).

We solve this convex program using first-order proximal methods with singular value shrinkage [33], implemented in `cvxpy`[1] with an SCS[2] solver. Step sizes are adapted with the Barzilai–Borwein [34] rule to accelerate convergence. The resulting $\hat{\Theta}$ provides continuous, psychometrically grounded labels for all 12 target dimensions.

## 3.3. Data Format and Structure

The RECRUITVIEW dataset comprises 2,011 multimodal samples, each representing a candidate's response to one of 76 interview questions. Each sample is structured to facilitate comprehensive multimodal analysis through three primary components:

- **Video**: High-resolution recordings stored in compressed MP4 format at 30 FPS. The dataset's average video duration is approximately 30 seconds.
- **Audio**: High-fidelity audio tracks extracted from videos (mono channel).
- **Transcript**: Verbatim speech-to-text transcriptions automatically generated using Whisper-large-v3[3] [35].

**Metadata and Annotations:** All annotations and metadata are organized in a structured JSON format. Each entry contains a unique identifier, video filename, interview question, quality indicators (video quality, duration category), user number and the 12 continuous target scores derived from the pairwise comparison protocol (see Appendix A.8 for a complete sample entry). This unified structure ensures seamless integration across modalities while maintaining data privacy and facilitating reproducible research workflows.

## 3.4. Task and Metrics

The primary task enabled by the RECRUITVIEW dataset is multimodal regression. Given a video clip of a candidate's response, the goal is to predict the 12 continuous scores corresponding to their personality traits and interview performance. Models leverage the modalities available from the data: visual (video frames), auditory (speech acoustics), and linguistic (transcribed text).

The 12 target variables for prediction are divided into the following two categories:

*Personality Traits Metrics:* These are based on the widely accepted Five-Factor Model of personality, with an additional overall score.

**1. Openness (O):** Measures imagination, creativity, and intellectual curiosity. Individuals high in openness are often inventive and enjoy new experiences.

**2. Conscientiousness (C):** Assesses self-discipline, organization, and goal-directed behavior. High conscientiousness is associated with being hardworking and reliable.

**3. Extraversion (E):** Reflects sociability, assertiveness, and emotional expressiveness. Extroverts tend to be outgoing and energized by social interaction.

**4. Agreeableness (A):** Indicates compassion, cooperativeness, and trustworthiness. Agreeable individuals are often helpful and empathetic.

**5. Neuroticism (N):** Pertains to emotional stability. Individuals high in neuroticism tend to experience negative emotions like anxiety and stress more frequently.

**6. Overall Personality:** A holistic assessment of the participant's perceived personality, derived from the combination of the Big Five traits.

*Performance Metrics:* These six metrics evaluate key competencies and behaviors exhibited during an interview response.

**7. Interview Score:** An overall score assessing the holistic quality of the participant's interview segment.

**8. Answer Score:** Evaluates the content of the response, including its relevance to the question, coherence, and structured thinking.

**9. Speaking Skills:** Assesses vocal characteristics such as clarity, pace, tone, and the avoidance of filler words.

**10. Confidence Score:** Measures the degree of self-assurance projected by the participant through both verbal and non-verbal cues (e.g., posture, eye contact, vocal tone).

**11. Facial Expression:** Quantifies the extent to which the participant uses facial expressions to convey emotion and engagement.

**12. Overall Performance:** A comprehensive evaluation of the candidate's performance in the clip, integrating all other performance factors.

## 3.5. Data Statistics

The RECRUITVIEW corpus comprises 2,011 video segments, sourced from over 300 unique participants responding to a bank of 76 curated interview questions. The dataset's foundation is a set of approximately 27,000 pairwise judgments provided by clinical psychologists. A key design characteristic is its "in-the-wild" data collection via a custom-built web-based platform. This methodology encouraged participation in naturalistic settings, resulting in significant variability in lighting conditions, background environments, and audio quality. This inherent diversity ensures high ecological validity, a critical feature for developing robust models that can generalize beyond controlled laboratory conditions. Figure 1 provides a summary of the video duration and quality categories, illustrating the distribution of these factors across the corpus.

To provide a more granular view of the dataset's temporal and linguistic composition, Figure 2 illustrates the distributions for video segment duration and the corresponding transcript word counts. The video durations (Figure 2, left) follow a right-skewed distribution, with a primary mode
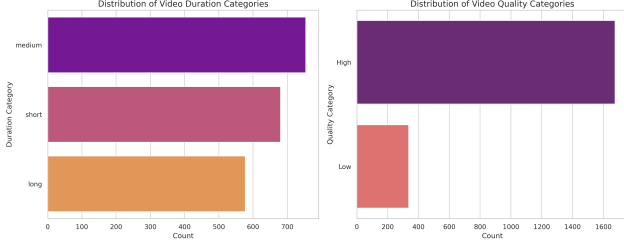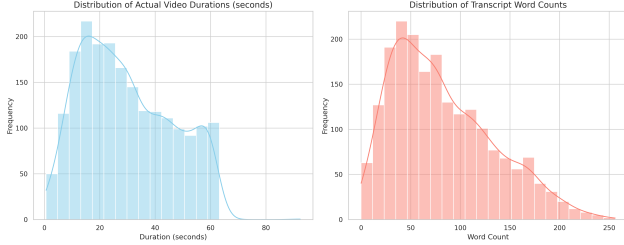
Figure 1. Data distribution by duration and video quality.



Figure 2. Data distribution by durations and transcript word count.



(a) Spearman's $\rho$ correlation matrix for personality trait metrics.

(b) Spearman's $\rho$ correlation matrix for performance metrics.

Figure 3. Spearman correlation between various metrics.

around 20-30 seconds and a secondary mode near 60 seconds, reflecting the natural variance in response length. The transcript word counts (Figure 2, right) are similarly distributed, with most responses falling between 40 and 115 words. This confirms that the dataset captures a wide range of response styles, from brief, concise answers to more detailed, elaborate ones.
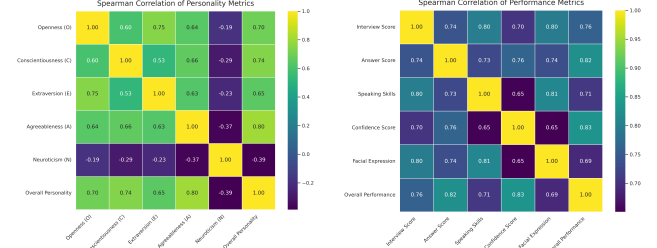
Detailed summary statistics for video durations and transcript word counts are provided in Appendix A.4, confirming the temporal and linguistic diversity of the dataset.

### 3.5.1. Correlation Analysis

To understand the relationships among the 12 target dimensions in RECRUITVIEW, we computed Spearman rank correlations across all 2,011 video clips.

**Personality Trait Metrics.** Figure 3a shows positive correlations between Overall Performance and Openness ($\rho = 0.70$), Conscientiousness ($\rho = 0.74$), Extraversion ($\rho = 0.65$), with Agreeableness strongest ($\rho = 0.80$); Neuroticism is negative ($\rho = -0.39$). This pattern aligns with theory: interpersonal warmth (Agreeableness) is most salient, while organization and intellectual engagement (Conscientiousness, Openness) also contribute.

**Performance Metrics.** Figure 3b shows moderate–strong positive correlations across all metrics. Confidence has the strongest association ($\rho = 0.83$), followed by Answer Score ($\rho = 0.82$) and Interview Score ($\rho = 0.76$). Facial Expressions and Speaking Skills are also substantial ($\rho = 0.69$, $\rho = 0.71$). Overall, stronger perceived personality aligns

with higher confidence, more expressive nonverbal behavior, and better-structured, well-delivered responses.

The complete $12 \times 12$ Spearman correlation matrix with all pairwise relationships is provided in Appendix A.5 for comprehensive reference.

### 3.5.2. Metrics Statistics

We analyzed the statistical properties of the 12 continuous target labels derived from the MNL model. The distributions are all normalized with near-zero means. A complete, detailed statistical analysis of all 12 metrics, including their implications, is provided in Appendix A.6. Most metrics exhibit significant leptokurtosis (heavy tails) and asymmetric skew. For instance, Speaking Skills ($\rho \approx -0.86$) and Overall Performance ($\rho \approx -0.75$) are negatively skewed, while Answer Score ($\rho \approx 0.35$) is positively skewed. This prevalence of outliers and non-normality strongly motivates our methodological choices: (1) the use of a robust **Huber loss** for training to mitigate the influence of extreme outliers, and (2) the prioritization of **rank-based correlation metrics** (e.g., Spearman's $\rho$) for evaluation, which are insensitive to this skew.

**Outlier Treatment via Soft Winsorization.** To address extreme outliers while preserving data structure, we apply mild soft winsorization. Values within $\pm 1.5\sigma$ pass through unchanged, while values beyond this threshold are smoothly compressed toward $\pm 3\sigma$ using tanh-based soft clipping: $\text{clip}(x) = \text{sign}(x) \cdot (\theta + s \cdot \tanh((|x| - \theta)/s))$ for $|x| > \theta$, where $\theta = 1.5$ and $s = 1.5$. This smooth transition prevents extreme values from dominating gradient updates while maintaining differentiability and rank ordering, improving convergence without discarding informative variance.

## 4. Methodology

### 4.1. Problem Formulation

We address the problem of predicting multiple continuous attributes from multimodal behavioral data. Given a video clip containing visual frames $\mathbf{V} \in \mathbb{R}^{T \times H \times W \times 3}$, audio

5

waveform $\mathbf{A} \in \mathbb{R}^L$, and transcript text $\mathbf{T}$, our goal is to predict a vector of target attributes $\mathbf{y} \in \mathbb{R}^K$ representing personality traits and performance scores. Formally, we learn a function $f : (\mathbf{V}, \mathbf{A}, \mathbf{T}) \rightarrow \mathbf{y}$ that captures the complex relationships between multimodal behavioral cues and target attributes.

Traditional approaches assume all representations reside in Euclidean space, employing linear transformations and standard neural operations. However, behavioral data exhibits diverse geometric properties: personality traits form hierarchical taxonomies (Big Five domains comprising specific facets), behavioral cues show directional relationships (facial expressions oriented in specific directions), and performance metrics often vary continuously. To capture this rich structure, we propose explicitly modeling representations across multiple geometric manifolds, each encoding different relational properties of the data.

## 4.2. CRMF Architecture Overview

The `CRMF` framework consists of six core components: (1) modality-specific encoders that extract representations from each input channel, (2) a pre-fusion module that performs early cross-modal integration, (3) manifold projection layers that map features to three geometric spaces, (4) geometry-specific expert networks that process each manifold representation, (5) an adaptive routing mechanism that learns optimal geometric combination strategies, and (6) a geometric fusion module that integrates multi-geometry representations for final prediction. Figure 4 illustrates the complete architecture.

The key insight behind `CRMF` is that different aspects of behavioral assessment benefit from different geometric inductive biases. Hyperbolic geometry naturally encodes hierarchical trait structures, spherical geometry captures directional behavioral patterns, and Euclidean geometry models continuous performance variations. By processing features through all three geometries and learning to combine them adaptively, `CRMF` can capture the full complexity of behavioral data. A detailed description of the `CRMF` framework's formulation and architecture is provided in Appendix D.

## 4.3. Multimodal Encoding

We employ pretrained encoders for each modality: DeBERTa-v3-base [36] for text, Wav2Vec2 [37]/HuBERT [38] for audio, and Video-MAE [39]/TimeSformer [40] for video. We fine-tune the last few layers of each encoder while keeping earlier layers frozen for parameter efficiency. For video, we apply a sophisticated temporal modeling pipeline consisting of BiLSTM, multi-head attention, and 1D convolution to capture rich temporal dynamics before pre-fusion. All modality encoders output representations with unified

dimension $d_{model} = 768$. Full details are provided in the Appendix D.1.

## 4.4. Pre-Fusion Module

The pre-fusion module performs early integration of multimodal features through cross-modal attention. We concatenate encoded features from all modalities and add learnable modality embeddings:

$$\mathbf{H}_{cat} = [\mathbf{H}_t; \mathbf{H}_a; \mathbf{H}_v] + \mathbf{E}_{mod} \tag{4}$$

where $\mathbf{E}_{mod} \in \mathbb{R}^{3 \times d}$ contains unique embeddings for text, audio, and video. A multi-layer transformer encoder processes the concatenated sequence, enabling rich cross-modal interactions. We employ learned attention pooling to obtain a fixed-dimensional clip-level representation $\mathbf{z}_{pre} \in \mathbb{R}^d$.

## 4.5. Manifold Projection and Geometric Experts

We project the fused representation $\mathbf{z}_{pre}$ onto three geometric manifolds using learned linear projections followed by geometry-specific mappings:

**Hyperbolic Space:** We use the Poincaré ball model $\mathbb{B}_c^d$ with curvature $c = 1.0$. Points are mapped via exponential map: $\exp_{\mathbf{0}}^c(\mathbf{v}) = \tanh(\sqrt{c}\|\mathbf{v}\|)\frac{\mathbf{v}}{\sqrt{c}\|\mathbf{v}\|}$, where $\mathbf{v} = \mathbf{W}_h\mathbf{z}_{pre}$.

**Spherical Space:** The unit sphere $\mathbb{S}^{d-1}$ is parameterized through $L_2$ normalization: $\mathbf{x}_s = \frac{\mathbf{W}_s\mathbf{z}_{pre}}{\|\mathbf{W}_s\mathbf{z}_{pre}\|+\epsilon}$.

**Euclidean Space:** Standard linear projection: $\mathbf{x}_e = \mathbf{W}_e\mathbf{z}_{pre}$.

Each manifold representation is processed by a specialized expert network designed to respect the underlying geometry. The hyperbolic expert uses Möbius transformations in gyrovector space, the spherical expert operates via tangent space mappings with exponential/logarithmic maps, and the Euclidean expert uses standard feed-forward layers. All experts have multiple layers and residual connections. Detailed formulations are available in the Appendix D.3.

## 4.6. Geometry-Aware Attention

To further refine expert outputs, we apply intra-manifold attention that respects geometric structure. For each geometry, we compute attention in its respective tangent space, which is Euclidean and enables standard multi-head attention operations. The attended representations are then mapped back to their respective manifolds.

## 4.7. Adaptive Token Routing

The router learns to weight expert outputs based on input characteristics. Given $\mathbf{z}_{pre}$, a lightweight MLP predicts routing weights:

$$\mathbf{r} = \text{softmax}(\mathbf{W}_r^{(2)}\sigma(\mathbf{W}_r^{(1)}\mathbf{z}_{pre} + \mathbf{b}_r^{(1)}) + \mathbf{b}_r^{(2)}) \tag{5}$$
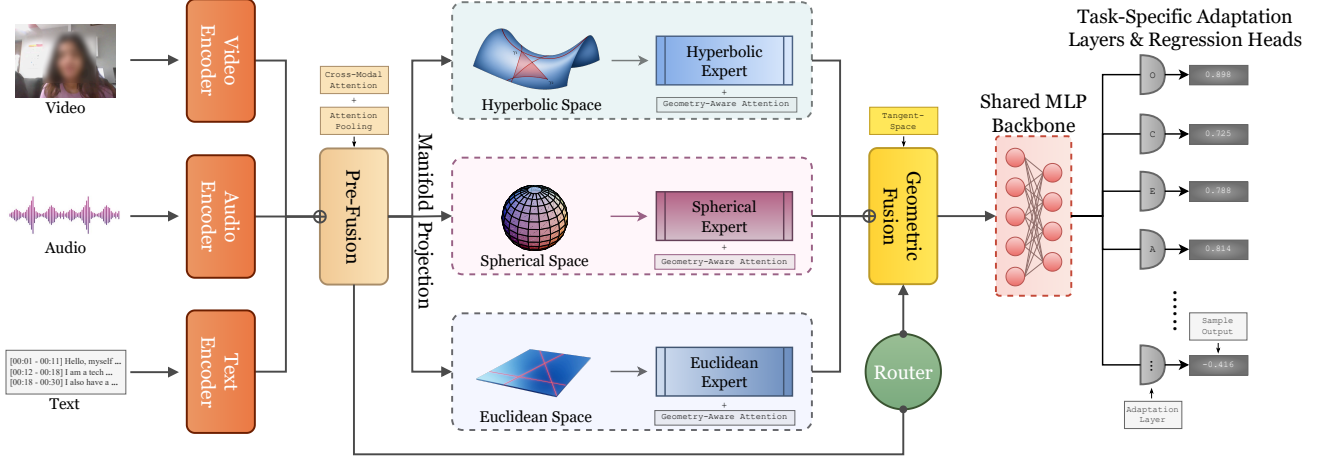
Figure 4. Overview of the CRMF architecture. Multimodal encoders extract features from video, audio, and text. Pre-fusion integrates modalities through cross-modal attention. The manifold projection layer maps features to hyperbolic, spherical, and Euclidean spaces. Geometry-specific experts process each manifold representation with intra-manifold attention. A learned router dynamically weights expert outputs. Finally, geometric fusion combines representations in a shared tangent space for multi-target prediction.

where $\mathbf{r} \in \Delta^{K-1}$ contains weights for the $K = 3$ experts. To encourage diverse geometry utilization, we apply entropy regularization: $\mathcal{L}_{entropy} = -\lambda_{ent} \sum_{i=1}^{K} r_i \log r_i$. A negative value encourages high entropy, promoting complementary geometric views.

### 4.8. Geometric Fusion

The fusion module combines expert outputs from different manifolds by first mapping all to a shared tangent space. For hyperbolic and spherical outputs, we use logarithmic maps; Euclidean output requires no conversion. The fusion operates via routing-weighted average:

$$\mathbf{z}_{fusion} = r_h \mathbf{v}_h + r_s \mathbf{v}_s + r_e \mathbf{v}_e \tag{6}$$

followed by a refinement network producing $\mathbf{z}_{refined}$. This strategy is equivalent to first-order Fréchet mean approximation on the product manifold.

### 4.9. Multi-Task Prediction Head

The prediction head maps the fused representation to target attributes using a shared MLP backbone followed by lightweight task-specific adaptation layers for each of the $K = 12$ targets. This parameter-efficient design balances expressiveness and efficiency.

### 4.10. Training Objective

Our training objective combines multiple loss components through adaptive balancing:

$$\mathcal{L}_{total} = \sum_{i=1}^{N} \beta_i \mathcal{L}_i \tag{7}$$

where components include Huber regression loss ($\delta = 1.0$), correlation boosting loss, covariance alignment loss, and auxiliary routing regularization losses. The weights $\beta_i$ are learned adaptively using inverse variance weighting combined with learned parameters. Full details are in the Appendix D.8.

## 5. Experimental Setup

### 5.1. Implementation Details

We train CRMF using AdamW [41] with component-specific learning rates. We use batch size 4 with gradient accumulation over 8 steps (effective batch size 32) and OneCycleLR scheduling with 15% warmup. Training runs for 30 epochs with early stopping on validation Spearman correlation (patience 5). Text is tokenized with max length 512, audio resampled to 16kHz, and video sampled at 16 FPS with 16 frames per clip resized to $224 \times 224$.

### 5.2. Baselines and Evaluation

We compare against three recent large multimodal models: MiniCPM-o 2.6 (8B) [1], VideoLLaMA2.1-AV (7B) [2], and Qwen2.5-Omni (7B) [3], all fine-tuned on our task. We evaluate using Spearman's $\rho$, Kendall's $\tau$-b, Concordance Index (C-Index), Pearson's $r$, and MSE. Metrics are computed per-target and macro-averaged for overall performance.

## 6. Results

### 6.1. Overall Performance

Table 1 presents aggregate results averaged across all 12 target attributes. CRMF variants substantially outperform

| Model | Spearman $\rho$ | Kendall $\tau$-b | C-index | Pearson $r$ |
|---|---|---|---|---|
| MiniCPM-o 2.6 (8B) | 0.5102 | 0.3541 | 0.6779 | 0.4935 |
| VideoLLaMA2.1-AV (7B) | 0.5002 | 0.3498 | 0.6778 | 0.4802 |
| Qwen2.5-Omni (7B) | 0.4882 | 0.3378 | 0.6658 | 0.4682 |
| CRMF (VMAE + w2v2) | **0.5682** | **0.4069** | **0.7183** | 0.5475 |
| CRMF (VMAE + HuB) | 0.5645 | 0.4020 | 0.7158 | **0.5481** |
| CRMF (TimeS + w2v2) | 0.5581 | 0.3980 | 0.7103 | 0.5394 |
| CRMF (TimeS + HuB) | 0.5664 | 0.4020 | 0.7148 | 0.5547 |

Table 1. Macro-averaged performance across all targets. CRMF variants consistently outperform baseline LMMs. Best results are bolded.

all baseline models across correlation and ranking metrics. The best CRMF configuration (VideoMAE + Wav2Vec2) achieves Spearman $\rho = 0.5682$, representing an 11.4% relative improvement over the strongest baseline (MiniCPM-o's 0.5102). Similar gains are observed for Kendall's $\tau$-b (14.9% improvement) and concordance index (6.0% improvement).

**Parameter Efficiency:** The performance gains are particularly remarkable given CRMF's parameter efficiency. Our framework (VideoMAE + Wav2Vec2 configuration) contains 408M parameters total, with 172M trainable during fine-tuning. In contrast, baseline LMMs fine-tune substantially more parameters: MiniCPM-o ($\sim$340M), VideoLLaMA2 ($\sim$300M), Qwen2.5-Omni ($\sim$320M). Despite training 40-50% fewer *trainable* parameters, CRMF achieves superior performance, demonstrating that task-specific geometric inductive biases provide more effective learning signals than simply leveraging larger pretrained models.

Comparing encoder choices, VideoMAE generally outperforms TimeSformer, suggesting that masked autoencoding provides better video representations for this task. For audio, Wav2Vec2 and HuBERT show comparable performance, with Wav2Vec2 having a slight edge on correlation metrics.

### 6.2. Per-Dimension Analysis

Per-trait personality assessment results (Table 8 in Appendix) show CRMF demonstrates strong performance across all Big Five dimensions. Openness shows the strongest CRMF performance ($\rho = 0.6384$), representing a 13.4% improvement over the best baseline. Conscientiousness, Extraversion, and Agreeableness exhibit moderate but consistent improvements (8-13% gains). Neuroticism presents the most challenging prediction task, though CRMF still improves upon baselines by 24.5%.

Per-dimension performance assessment results (Table 9 in Appendix) detail results for six performance-related attributes. Interview Score and Answer Score show the strongest absolute performance ($\rho > 0.62$ and $\rho > 0.59$), with 9-12% improvements over baselines. Speaking Skills

| Component | Variant | $\rho$ | $\tau$-b | C-idx |
|---|---|---|---|---|
| *Full CRMF Model* | | 0.5682 | 0.4069 | 0.7183 |
| **Fusion Only** | Simple Concatenation | 0.4441 | 0.2903 | 0.6151 |
| | Weighted Average | 0.4664 | 0.3078 | 0.6239 |
| **Geometry** | Hyperbolic Only | 0.5080 | 0.3611 | 0.6980 |
| | Spherical Only | 0.5338 | 0.3808 | 0.7054 |
| | Euclidean Only | 0.5284 | 0.3753 | 0.7001 |
| **Routing** | No Router (Uniform) | 0.5209 | 0.3688 | 0.6994 |
| **Modality** | Video Only (VMAE) | 0.4516 | 0.2974 | 0.6521 |
| | Audio Only (Wav2Vec2) | 0.3792 | 0.2461 | 0.6245 |
| | Text Only (DeBERTa) | 0.4247 | 0.2806 | 0.6429 |

Table 2. Key ablation study results demonstrating the contribution of CRMF components. All experiments use Video-MAE+Wav2Vec2.

and Confidence Score achieve moderate but consistent improvements (10-16% gains). Overall Performance benefits most from CRMF ($\rho = 0.6521$), with 9.1% improvement.

### 6.3. Ablation Studies

Table 2 presents key ablation results. Using only fusion (no CRMF framework), such as simple concatenation ($\rho = 0.4441$) or weighted averaging ($\rho = 0.4664$), causes severe degradation (21.8% and 17.9% drops), demonstrating that naive fusion strategies fail to capture complex multimodal relationships.

Using only a single geometric space consistently underperforms: Hyperbolic-only achieves $\rho = 0.5080$ (10.6% drop), Spherical-only $\rho = 0.5338$ (6.1% drop), and Euclidean-only $\rho = 0.5284$ (7.0% drop). No single geometry matches the full model, confirming that different geometric spaces capture complementary information.

Removing the router and using uniform weights ($\rho = 0.5209$) causes substantial degradation (8.3% drop), confirming that adaptive weighting based on input characteristics is crucial.

Single modality analysis reveals video provides the strongest unimodal signal ($\rho = 0.4516$), followed by text ($\rho = 0.4247$) and audio ($\rho = 0.3792$). Crucially, even the best unimodal result is substantially lower than any multimodal configuration. The leap from video-only to full CRMF represents a 25.8% improvement, underscoring that behavioral traits are expressed through complex interplay of cues across modalities.

Additional ablation results for two-geometry combinations, pre-fusion variants, prediction head architectures, and loss functions are provided in the Appendix (Table 10).

## 7. Conclusion

We introduced RECRUITVIEW, a multimodal corpus for personality and interview analysis with continuous, psychometrically grounded labels derived from pairwise judgments. Building on this resource, we proposed CRMF, a geometry-aware regression framework that fuses audio, video, and text using manifold-specific attention and adaptive routing. On RECRUITVIEW, CRMF surpasses strong multimodal baselines, raising macro Spearman correlation to 0.568 and concordance index to 0.718, while using fewer trainable parameters. Ablations validate the benefits of multi-geometry fusion and routing, and show clear gaps between multimodal and unimodal variants. Limitations include moderate dataset scale, short clips, potential residual annotator bias and label noise despite calibration, and limited demographic diversity, which may constrain external validity. Future work will broaden populations and conditions, extend to longer multi-turn interactions, and integrate stronger self-supervised priors, target-wise manifold selection, and causal analyses, alongside real-time inference and human-in-the-loop calibration.

## Data and Code Availability

The RECRUITVIEW dataset and the CRMF implementation are publicly available. The dataset is hosted on Hugging Face at https://huggingface.co/datasets/AI4A-lab/RecruitView and mirrored on GitHub at https://github.com/AI4A-lab/RecruitView. The CRMF framework code is available at https://github.com/AI4A-lab/CRMF.

## References

[1] S. Hu, Y. Tu, X. Han, C. He, G. Cui, X. Long, Z. Zheng, Y. Fang, Y. Huang, W. Zhao, *et al.*, "Minicpm: Unveiling the potential of small language models with scalable training strategies," *arXiv preprint arXiv:2404.06395*, 2024. 1, 2, 7

[2] Z. Cheng, S. Leng, H. Zhang, Y. Xin, X. Li, G. Chen, Y. Zhu, W. Zhang, Z. Luo, D. Zhao, *et al.*, "Videollama 2: Advancing spatial-temporal modeling and audio understanding in video-llms," *arXiv preprint arXiv:2406.07476*, 2024. 1, 2, 7

[3] J. Xu, Z. Guo, J. He, H. Hu, T. He, S. Bai, K. Chen, J. Wang, Y. Fan, K. Dang, *et al.*, "Qwen2.5-omni technical report," *arXiv preprint arXiv:2503.20215*, 2025. 1, 2, 7

[4] D. P. McAdams, "The five-factor model in personality: A critical appraisal," *Journal of personality*, vol. 60, no. 2, pp. 329–361, 1992. 1

[5] V. Ponce-López, B. Chen, M. Oliu, C. Corneanu, A. Clapés, I. Guyon, X. Baró, H. J. Escalante, and S. Escalera, "Chalearn lap 2016: First round challenge on first impressions-dataset and results," in *European conference on computer vision*, pp. 400–418, Springer, 2016. 2, 3

[6] S. Park, H. S. Shim, M. Chatterjee, K. Sagae, and L.-P. Morency, "Computational analysis of persuasiveness in social multimedia: A novel dataset and multimodal prediction approach," in *Proceedings of the 16th International Conference on Multimodal Interaction*, ICMI '14, (New York, NY, USA), p. 50–57, Association for Computing Machinery, 2014. 2

[7] J.-I. Biel and D. Gatica-Perez, "Facetube: Predicting personality from facial expressions of emotion in online conversational video," in *Proceedings of the 13th International Conference on Multimodal Interaction (ICMI)*, pp. 53–56, ACM, 2011. 2

[8] Y. Song, F. Yang, S. Huang, and L. Chen, "Interview2personality: A multimodal dataset for predicting personality traits from job interviews," *IEEE Transactions on Affective Computing*, 2020. 2

[9] T. Baltrušaitis, C. Ahuja, and L.-P. Morency, "Multimodal machine learning: A survey and taxonomy," *IEEE transactions on pattern analysis and machine intelligence*, vol. 41, no. 2, pp. 423–443, 2018. 2

[10] Y.-H. H. Tsai, S. Bai, P. P. Liang, J. Z. Kolter, L.-P. Morency, and R. Salakhutdinov, "Multimodal transformer for unaligned multimodal language sequences," in *Proceedings of the conference. Association for computational linguistics. Meeting*, vol. 2019, p. 6558, 2019. 2

[11] A. Zadeh, M. Chen, S. Poria, E. Cambria, and L.-P. Morency, "Tensor fusion network for multimodal sentiment analysis," *arXiv preprint arXiv:1707.07250*, 2017. 2

[12] O. Ganea, G. Bécigneul, and T. Hofmann, "Hyperbolic neural networks," *Advances in neural information processing systems*, vol. 31, 2018. 2

[13] I. Chami, Z. Ying, C. Ré, and J. Leskovec, "Hyperbolic graph convolutional neural networks," *Advances in neural information processing systems*, vol. 32, 2019. 2

[14] T. S. Cohen, M. Geiger, J. Köhler, and M. Welling, "Spherical cnns," *arXiv preprint arXiv:1801.10130*, 2018. 2

[15] C. Esteves, C. Allen-Blanchette, A. Makadia, and K. Daniilidis, "Learning so (3) equivariant representations with spherical cnns," in *Proceedings of the european conference on computer vision (ECCV)*, pp. 52–68, 2018. 2

[16] A. Gu, F. Sala, B. Gunel, and C. Ré, "Learning mixed-curvature representations in product spaces," in *International conference on learning representations*, 2018. 2

[17] O. Skopek, O.-E. Ganea, and G. Bécigneul, "Mixed-curvature variational autoencoders," *arXiv preprint arXiv:1911.08411*, 2019. 2

[18] D. Brooks, O. Schwander, F. Barbaresco, J.-Y. Schneider, and M. Cord, "Riemannian batch normalization for spd neural networks," *Advances in Neural Information Processing Systems*, vol. 32, 2019. 2

[19] A. Lou, D. Lim, I. Katsman, L. Huang, Q. Jiang, S. N. Lim, and C. M. De Sa, "Neural manifold ordinary differential equations," *Advances in Neural Information Processing Systems*, vol. 33, pp. 17548–17558, 2020. 2

[20] N. Shazeer, A. Mirhoseini, K. Maziarz, A. Davis, Q. Le, G. Hinton, and J. Dean, "Outrageously large neural networks: The sparsely-gated mixture-of-experts layer," *arXiv preprint arXiv:1701.06538*, 2017. 2

[21] W. Fedus, B. Zoph, and N. Shazeer, "Switch transformers: Scaling to trillion parameter models with simple and efficient sparsity," *Journal of Machine Learning Research*, vol. 23, no. 120, pp. 1–39, 2022. 2

[22] B. Mustafa, C. Riquelme, J. Puigcerver, R. Jenatton, and N. Houlsby, "Multimodal contrastive learning with limoe: the language-image mixture of experts," *Advances in Neural Information Processing Systems*, vol. 35, pp. 9564–9576, 2022. 2

[23] J. Xin, S. Yun, J. Peng, I. Choi, J. L. Ballard, T. Chen, and Q. Long, "I2moe: Interpretable multimodal interaction-aware mixture-of-experts," *arXiv preprint arXiv:2505.19190*, 2025. 2

[24] B. Chen, S. Escalera, I. Guyon, V. Ponce-López, N. Shah, and M. Oliu Simón, "Overcoming calibration problems in pattern labeling with pairwise ratings: application to personality traits," in *European Conference on Computer Vision*, pp. 419–432, Springer, 2016. 3

[25] L. L. Thurstone, "A law of comparative judgment," *Psychological Review*, vol. 34, no. 4, pp. 273–286, 1927. 3

[26] M. E. Glickman and A. C. Jones, "Rating the chess rating system," *Chance*, vol. 12, no. 2, pp. 21–28, 1999. 3

[27] R. A. Bradley and M. E. Terry, "Rank analysis of incomplete block designs. I. The method of paired comparisons," *Biometrika*, vol. 39, no. 3/4, pp. 324–345, 1952. 3

[28] R. D. Luce, *Individual choice behavior: A theoretical analysis*. Wiley, 1959. 3

[29] R. Herbrich, T. Minka, and T. Graepel, "TrueSkill: A Bayesian Skill Rating System," in *Advances in Neural Information Processing Systems 19 (NIPS)*, pp. 569–576, 2007. 3

[30] M. E. Glickman, "Parameter estimation in large dynamic paired comparison experiments," *Applied Statistics*, vol. 48, no. 3, pp. 377–394, 1999. 3

[31] S. Negahban, S. Oh, K. K. Thekumparampil, and J. Xu, "Learning from Comparisons and Choices," *Journal of Machine Learning Research*, vol. 19, pp. 1–95, 2018. 3

[32] M. Fazel, *Matrix rank minimization with applications*. PhD thesis, Stanford University, 2002. 3

[33] J.-F. Cai, E. J. Candès, and Z. Shen, "A singular value thresholding algorithm for matrix completion," *SIAM Journal on Optimization*, vol. 20, no. 4, pp. 1956–1982, 2010. 4

[34] J. Barzilai and J. M. Borwein, "Two-point step size gradient methods," *IMA journal of numerical analysis*, vol. 8, no. 1, pp. 141–148, 1988. 4

[35] A. Radford, J. W. Kim, T. Xu, G. Brockman, C. McLeavey, and I. Sutskever, "Robust speech recognition via large-scale weak supervision," in *International conference on machine learning*, pp. 28492–28518, PMLR, 2023. 4

[36] P. He, J. Gao, and W. Chen, "Debertav3: Improving deberta using electra-style pre-training with gradient-disentangled embedding sharing," *arXiv preprint arXiv:2111.09543*, 2021. 6, 17

[37] A. Baevski, Y. Zhou, A. Mohamed, and M. Auli, "wav2vec 2.0: A framework for self-supervised learning of speech representations," *Advances in neural information processing systems*, vol. 33, pp. 12449–12460, 2020. 6, 17

[38] W.-N. Hsu, B. Bolte, Y.-H. H. Tsai, K. Lakhotia, R. Salakhutdinov, and A. Mohamed, "Hubert: Self-supervised speech representation learning by masked prediction of hidden units," *IEEE/ACM transactions on audio, speech, and language processing*, vol. 29, pp. 3451–3460, 2021. 6, 17

[39] Z. Tong, Y. Song, J. Wang, and L. Wang, "Videomae: Masked autoencoders are data-efficient learners for self-supervised video pre-training," *Advances in neural information processing systems*, vol. 35, pp. 10078–10093, 2022. 6, 18

[40] G. Bertasius, H. Wang, and L. Torresani, "Is space-time attention all you need for video understanding?," in *Proceedings of the International Conference on Machine Learning (ICML)*, July 2021. 6, 18

[41] I. Loshchilov and F. Hutter, "Decoupled weight decay regularization," *arXiv preprint arXiv:1711.05101*, 2017. 7

[42] A. Ungar, *A gyrovector space approach to hyperbolic geometry*. Springer Nature, 2022. 19

## A. RECRUITVIEW Dataset

### A.1. Questions

The full list of the 76 unique interview questions used as prompts in the data collection is provided in Table 3. These questions were curated from professional open-source resources, networking platforms, and consultations to elicit responses rich in both professional content and personality indicators.

### A.2. Collection and Labeling Framework

To ensure standardized data acquisition and annotation, we developed two custom web-based platforms. The first, `QAVideoShare`[4], is an online interview platform designed to collect video responses. Participants were presented with questions and recorded unscripted answers directly through the browser interface, ensuring uniform question presentation and automated video storage. The participant's workflow, from authentication to recording, is shown in Figure 5.

The second platform, `QA-Labeler`[5], was developed for data labeling and evaluation. This tool allowed clinical psychologists (annotators) to view recorded videos and provide comparative assessments across various behavioral and performance criteria. The comparative judgment interface, featuring a side-by-side player and scoring form, is detailed in Figure 6. Both platforms support browser-based, multi-user operation, enabling a scalable and consistent data processing pipeline.

### A.3. Labeling Comparison

To ensure that the conversion of pairwise judgments into continuous rankings was both consistent and interpretable, we evaluated five independent label-conversion frameworks: **Glicko-2**, **TrueSkill**, **Full-Rank MNL**, **MNL-with-Ties**, and the **Nuclear-Norm-Regularized MNL**. Each method produced a scalar ranking score representing the latent position of each video across all pairwise comparisons. All models were trained on identical data.

#### A.3.1. Ground-Truth Construction

We adopt a leave-one-out consensus evaluation. When evaluating a given model, its output is compared against a ground-truth defined as the mean of the standardized (Z-scored) scores from all *other* models. This avoids self-evaluation, ensures symmetric treatment of all frameworks, and controls for scale or range mismatches. For models such as Nuclear-Norm MNL which inherently output normalized scores, further standardization was not applied.

---

#### A.3.2. Evaluation Results

Evaluation was performed using **Spearman's** $\rho$, **Kendall's** $\tau$, **RMSE**, **MAE**, and **Precision/Recall@10%**, computed on continuous ranking outputs. These metrics were chosen for their compatibility with ordinal data, as our objective is to evaluate *relative ordering* fidelity rather than categorical correctness. Figure 7 shows the Spearman rank-correlation structure across models. Here, higher correlation is desirable because label-conversion methods are not supposed to invent disagreement; divergence between methods would indicate instability or method-specific distortion rather than genuine latent behavioral differences. Figure 8 further reinforces this observation, showing all five models plotted against the consensus reference for direct visual comparison.

Although all frameworks produced broadly compatible rankings, **Full-Rank MNL** achieved slightly higher peak correlations on isolated traits, while the **Nuclear-Norm MNL** exhibited greater overall stability with low variance across random drop-model trials ($\rho = 0.905 \pm 0.04$). The low-rank constraint enforces smoother coupling across correlated personality traits, yielding more stable global rankings; this behavior is further reflected in the robustness summary (Table 4).

#### A.3.3. Secondary Verification and Qualitative Assessment

A secondary verification was conducted through a ratio-based test: a manually selected subset of pairwise comparisons was converted into empirical win–loss ratios, which serve as a local ordinal reference. The Nuclear-Norm MNL produced the closest match, accurately preserving both relative order and proportional differences. A small leaderboard test confirmed that local chains (A > B > C) remained globally consistent (A > C) and aligned with human expectations. In qualitative inspection, videos ranked higher by this model displayed clearer articulation, stronger confidence, and more natural expressiveness.

Accordingly, we adopt the Nuclear-Norm formulation as the final label-conversion framework for RECRUITVIEW. Its low-rank structure offered smoother scaling across correlated targets, and its predictions were the most consistent with manual verification.

### A.4. Detailed Data Statistics

Table 5 presents comprehensive summary statistics for the 2,011 video segments in RECRUITVIEW. The clips have a mean duration of 29.66 seconds ($\sigma = 16.40$), with a minimum of 0.60 seconds and a maximum of 92.34 seconds. This temporal range ensures models are exposed to both brief "thin-slice" judgments and longer-form analyses. The transcripts are similarly diverse, with a mean word count of 81.90 ($\sigma = 51.15$) and a maximum of 266 words, providing

| Sr. No. | Question | | Sr. No. | Question |
|---|---|---|---|---|
| 1. | Introduce yourself. | | 41. | What are the three things that are most important for you in a job? |
| 2. | What are your greatest strengths and weaknesses? | | 42. | How did you handle disagreements? |
| 3. | How do you handle changes or unexpected situations in the workplace? | | 43. | Tell me about a time where you experienced difficulty while working on a project. How did you handle it? |
| 4. | What is your biggest achievement so far? | | 44. | What makes you happy? |
| 5. | Tell me about a time when you went above and beyond the call of duty to achieve a goal or deliver results. | | 45. | Can you give an example of a situation where you mentored a junior colleague, helping them grow professionally and personally? |
| 6. | Give me an example of your creativity. | | 46. | What are you passionate about? |
| 7. | How do you work under pressure? Can you handle the pressure? | | 47. | What motivates you to perform at your best in the workplace? |
| 8. | If you won a Rs.10-crore lottery, would you still work? | | 48. | Describe a time when you proactively sought out opportunities to develop new skills or knowledge relevant to your role. |
| 9. | What motivates you? | | 49. | Can you give an example of a situation where you leveraged technology or automation to streamline a process and increase efficiency? |
| 10. | Can you give an example of a time when you successfully implemented a solution to improve a process or procedure? | | 50. | Share a story of a project where you collaborated with a cross-functional team to deliver exceptional results. |
| 11. | Tell me about a time when you had to step into a role outside of your expertise to support the team's objectives. | | 51. | Describe a situation where you had to adapt to a change in the work environment. |
| 12. | How do you respond to change? | | 52. | What are you most proud of? |
| 13. | What was the toughest decision you ever had to make? | | 53. | What do you think is an ideal work environment? |
| 14. | What is your greatest fear? | | 54. | Tell me about a time you initiated or led that had a positive impact on your team or organization. |
| 15. | Describe a project where you took the lead in implementing a new strategy or process, driving positive change within your team or organization. | | 55. | Can you give an example of a time when you had to resolve a disagreement or misunderstanding within a team? |
| 16. | Describe a situation where you identified a problem before it became significant. What steps did you take to address it? | | 56. | How do you deal with criticism? |
| 17. | How would you rate yourself on a scale of 1 to 10? | | 57. | Tell me about a time when you failed to meet a goal or objective. How did you handle it? |
| 18. | How do you handle stress and anxiety? | | 58. | What has been your greatest failure? |
| 19. | Tell me about a time when you were not satisfied with your performance. | | 59. | Tell me about a time when you had to resolve a conflict with a coworker or team member. |
| 20. | Can you give an example of a time when you successfully managed multiple competing priorities? | | 60. | Share a story of a project where you led the team in developing and implementing a solution that resulted in significant cost savings or revenue growth. |
| 21. | Where do you see yourself in the next 5 years? | | 61. | Tell me about a time you had to work on a project outside of your comfort zone. How did you handle it? |
| 22. | Why should a company hire you? | | 62. | Are you an organized person? |
| 23. | Are you reliable or can I trust you with responsibilities? | | 63. | What do you always regret, or do you have any regrets? |
| 24. | What makes you angry? | | 64. | Share a story of how you took ownership of a project that was struggling and turned it into a success through your initiative. |
| 25. | Can you give an example of a time when you had to persuade others to adopt your ideas or proposals? | | 65. | Can you give an example of a situation where you successfully motivated a disengaged team member to contribute effectively to a project? |
| 26. | Can you give an example of a time when you had to take the initiative to solve a problem without being asked? | | 66. | How do you learn new skills? |
| 27. | Are you open to taking risks or do you like experimenting? | | 67. | Describe a situation where you had to prioritize tasks under tight deadlines. |
| 28. | What is your dream company like? | | 68. | How quickly do you adapt to new technology? |
| 29. | Do you have a good work ethic? | | 69. | Can you give an example of a time when you facilitated a productive team meeting or discussion? |
| 30. | Tell me about a time when you recognized and capitalized on the unique strengths of individual team members to achieve a common goal. | | 70. | Describe a project where you collaborated with stakeholders to define the problem and develop a solution that met everyone's needs. |
| 31. | Tell me about a time when you successfully resolved a long-standing issue that had been impeding progress within your team or organization. | | 71. | What is your dream job like? |
| 32. | How do you improve your knowledge? | | 72. | What are your weaknesses? |
| 33. | What are your hobbies? | | 73. | Tell me about a time when you proposed an innovative idea that significantly improved team efficiency or productivity. |
| 34. | Can you give an example of a time when you led by example to promote a positive work culture or values? | | 74. | Can you give an example of a time when you coached or mentored a colleague to help them achieve their goals? |
| 35. | Describe a project where you encouraged open communication and feedback among team members, leading to improved collaboration and outcomes. | | 75. | Share a story of a time when you rallied your team during a crisis, fostering resilience and determination. |
| 36. | Is there anything that makes you different from other candidates? | | 76. | Tell me about yourself. |
| 37. | Can you describe your time management skills? | | | |
| 38. | Can you describe a situation where you had to overcome a significant challenge in a team setting? | | | |

Table 3. The 76 curated interview questions used as prompts in the RECRUITVIEW dataset.
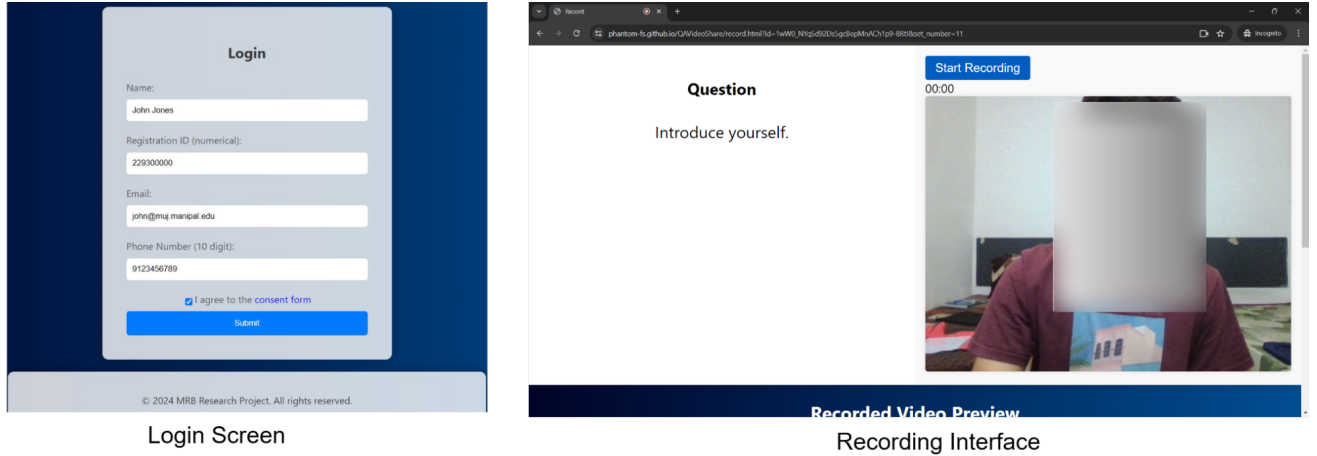


Figure 5. The participant-facing `QAVideoShare` data collection platform. (Left) The secure login and consent portal. (Right) The primary video recording interface where participants view the prompt and record their response.
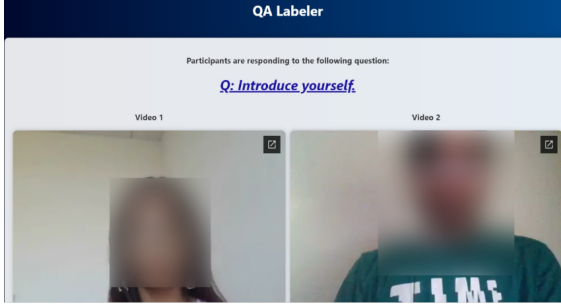
| Model | Avg. $\rho$ | Std. Dev. | Trials |
|---|---|---|---|
| MNL (Full Rank) | 0.910 | 0.044 | 13 |
| MNL (with Ties) | 0.903 | 0.042 | 17 |
| MNL (Nuclear Norm) | 0.905 | 0.040 | 19 |
| Glicko-2 | 0.860 | 0.049 | 16 |
| TrueSkill | 0.776 | 0.053 | 15 |

Table 4. Robustness check across 20 random drop-model trials.

| Statistic | Duration (seconds) | Word Count |
|---|---|---|
| count | 2011.00 | 2011.00 |
| mean | 29.66 | 81.90 |
| std | 16.40 | 51.15 |
| min | 0.60 | 0.00 |
| 25% | 15.83 | 41.00 |
| 50% | 27.27 | 72.00 |
| 75% | 42.63 | 115.00 |
| max | 92.34 | 256.00 |

Table 5. Statistics for Engineered Features in RECRUITVIEW. The table shows distribution statistics for video duration (in seconds) and transcript word count across all 2,011 clips.

a rich linguistic substrate for multimodal analysis. The median (50th percentile) values for duration (27.27s) and word count (72.00) closely track their respective means, confirming the well-behaved nature of these distributions.

Figure 6. The evaluator-facing `QA-Labeler` annotation platform. (Left) The side-by-side video playback module for comparative judgment. (Right) The corresponding scoring form where evaluators provide choice ratings on behavioral traits.
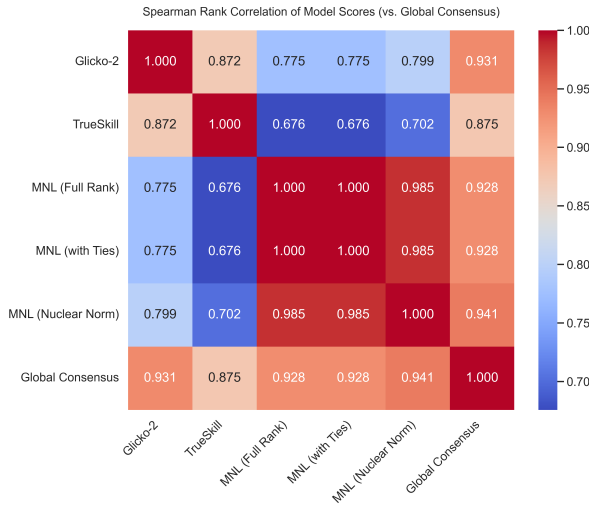


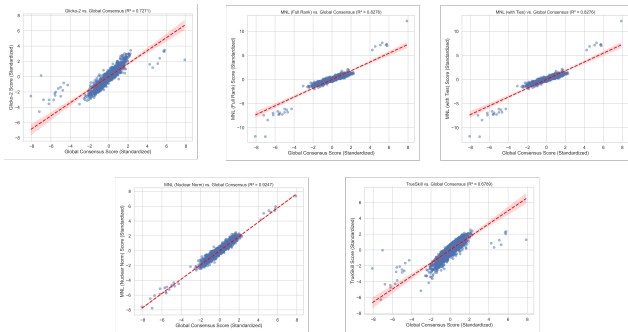Figure 7. Correlation matrix among the five label-conversion frameworks.



Figure 8. Scatter plots of predicted vs. ground-truth rankings for all five label-conversion frameworks. Each subplot corresponds to one model.

## A.5. Complete Correlation Matrix

Table 6 presents the complete Spearman correlation matrix across all 12 target dimensions in RECRUITVIEW. This

comprehensive view consolidates the patterns observed in Figure 9, revealing the full structure of dependencies among the Big Five personality traits (O, C, E, A, N), Overall Personality, and the six interview performance metrics (Interview Score, Answer Score, Speaking Skills, Confidence Score, Facial Expression, and Overall Performance). The matrix exhibits several key characteristics: strong positive correlations within the personality cluster (upper-left block) and performance cluster (bottom-right block), moderate positive correlations in the cross-domain blocks, and consistent negative correlations involving Neuroticism across all dimensions. The cross-correlation block (bottom-left) shows intuitive patterns:

- *Extraversion* is positively correlated with *Speaking Skills* ($\rho = 0.71$) and *Facial Expression* ($\rho = 0.71$), suggesting that outgoing individuals are perceived as more expressive and articulate.
- *Conscientiousness* shows a clear positive relationship with *Answer Score* ($\rho = 0.70$), aligning with the expectation that diligent individuals provide higher-quality responses.
- *Neuroticism* demonstrates a consistent negative correlation across all performance metrics, most notably with *Confidence Score* ($\rho = -0.37$) and *Overall Performance* ($\rho = -0.36$).

These structured dependencies highlight the interconnected nature of personality perception and observable interview behaviors, providing insights into which combinations of traits and performance indicators are most strongly linked in evaluative contexts.

## A.6. Metrics Statistics

Table 7 and Figure 10 jointly characterize the empirical behavior of all twelve targets in RECRUITVIEW. First, the *means* sit essentially at zero for every dimension (see "mean" row), confirming that the normalization pipeline yields centered targets and ensuring that absolute intercepts are not informative. The *medians* closely track the means

13

| Metrics | O | C | E | A | N | Pers. | Int. | Ans. | Spk. | Conf. | Fac. | Perf. |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Openness (**O**) | 1.00 | 0.60 | 0.75 | 0.64 | -0.19 | 0.70 | 0.77 | 0.68 | 0.74 | 0.67 | 0.75 | 0.72 |
| Conscientiousness (**C**) | 0.60 | 1.00 | 0.53 | 0.66 | -0.29 | 0.74 | 0.59 | 0.70 | 0.59 | 0.65 | 0.59 | 0.70 |
| Extraversion (**E**) | 0.75 | 0.53 | 1.00 | 0.63 | -0.23 | 0.65 | 0.71 | 0.67 | 0.71 | 0.68 | 0.71 | 0.66 |
| Agreeableness (**A**) | 0.64 | 0.66 | 0.63 | 1.00 | -0.37 | 0.80 | 0.69 | 0.78 | 0.70 | 0.73 | 0.70 | 0.77 |
| Neuroticism (**N**) | -0.19 | -0.29 | -0.23 | -0.37 | 1.00 | -0.39 | -0.27 | -0.37 | -0.28 | -0.37 | -0.26 | -0.36 |
| Overall Personality (**Pers.**) | 0.70 | 0.74 | 0.65 | 0.80 | -0.39 | 1.00 | 0.73 | 0.79 | 0.69 | 0.80 | 0.71 | 0.84 |
| Interview Score (**Int.**) | 0.77 | 0.59 | 0.71 | 0.69 | -0.27 | 0.73 | 1.00 | 0.74 | 0.80 | 0.70 | 0.80 | 0.76 |
| Answer Score (**Ans.**) | 0.68 | 0.70 | 0.67 | 0.78 | -0.37 | 0.79 | 0.74 | 1.00 | 0.73 | 0.76 | 0.74 | 0.82 |
| Speaking Skills (**Spk.**) | 0.74 | 0.59 | 0.71 | 0.70 | -0.28 | 0.69 | 0.80 | 0.73 | 1.00 | 0.65 | 0.81 | 0.71 |
| Confidence Score (**Conf.**) | 0.67 | 0.65 | 0.68 | 0.73 | -0.37 | 0.80 | 0.70 | 0.76 | 0.65 | 1.00 | 0.65 | 0.83 |
| Facial Expression (**Fac.**) | 0.75 | 0.59 | 0.71 | 0.70 | -0.26 | 0.71 | 0.80 | 0.74 | 0.81 | 0.65 | 1.00 | 0.69 |
| Overall Performance (**Perf.**) | 0.72 | 0.70 | 0.66 | 0.77 | -0.36 | 0.84 | 0.76 | 0.82 | 0.71 | 0.83 | 0.69 | 1.00 |

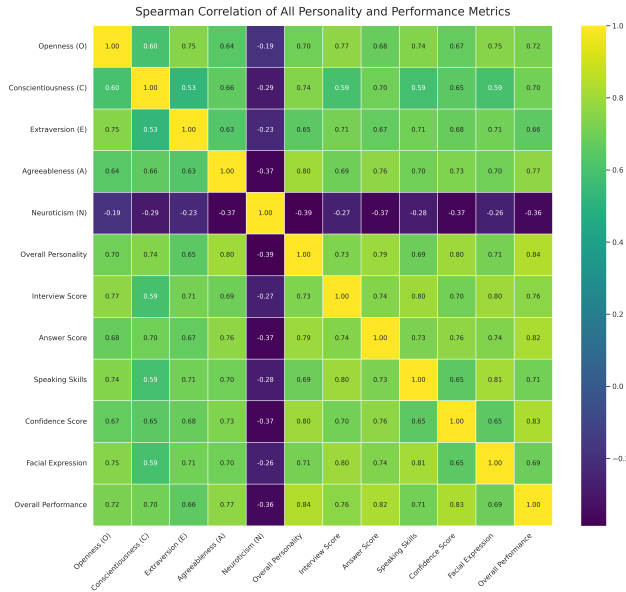Table 6. Complete Spearman Correlation Matrix for all 12 metrics in RECRUITVIEW.



Figure 9. Spearman's $\rho$ correlation matrix for all 12 metrics.

(50% row in Table 7), and the modal mass of each histogram is concentrated around the origin (Figure 10), indicating a near-symmetric *core* for most variables.

**Dispersion and dynamic range.** Standard deviations cluster in the $[0.88, 1.28]$ interval for 11/12 metrics, with **Neuroticism (N)** exhibiting a much tighter spread ($\text{std} = 0.49$). This implies that N is intrinsically less variable across our population relative to other psychological or performance attributes. Conversely, Interview-adjacent outcomes (Int., Ans., Spk., Conf., Fac., Perf.) show broadly comparable dispersion ($\approx 1.1$–$1.28$), desirable for multi-task optimization with shared heads. The *extrema* reveal long tails for several metrics (e.g., Ans.: $\min = -10.20$; O: $\max = 9.34$), which are far beyond $\pm 3\sigma$ and thus consti-

tute influential observations for any squared-loss estimator.
**Asymmetry (skewness).** Skewness in Table 7 uncovers systematic asymmetries:

- *Negative skew* for **C**, **Spk**, **Conf**, and **Perf** $(-0.57, -0.86, -0.64, -0.75)$ indicates heavier left tails and a right-shifted bulk. Practically, a larger fraction of samples achieve above-average performance on speaking, confidence, and overall performance, with relatively fewer but more extreme low outliers.

- *Positive skew* for **A**, **Pers.**, **Int.**, and **Fac.** $(0.40$–$0.66)$ suggests the opposite: mass slightly left of zero with occasional high outliers. **Openness** and **Extraversion** show mild asymmetry $(0.03$ and $-0.22)$, while **Neuroticism** is modestly negative $(-0.25)$, again consistent with its compressed variance.

These asymmetries imply that symmetric error models may under- or over-penalize different tails across tasks; model selection should therefore consider robust losses and rank-based metrics.

**Tail heaviness (kurtosis).** All targets except **Neuroticism** show pronounced leptokurtosis (excess kurtosis $\approx$ 8.8–13.4), confirming heavy tails and a high concentration near the center (Table 7). **Neuroticism** (1.14) is notably closer to mesokurtic behavior relative to the other metrics. Combined with the extreme minima/maxima, this indicates that a small subset of clips carry disproportionately informative deviations—a regime where (i) Huber/quantile losses and (ii) clipping or winsorization, materially improve stability and interpretability.

**Quartiles and central mass.** Interquartile ranges are tightly packed around zero (25%–75% roughly $\pm 0.55$–$0.62$ for most metrics), reinforcing that the majority of ratings occupy a narrow band. The practical upshot is twofold: (a) small absolute errors around the origin correspond to meaningful rank changes, and (b) evaluation should prioritize *monotonicity* (*Spearman $\rho$*, *Kendall $\tau$*, or concordance

index) in addition to pointwise deviations.

## A.7. Data Splits

We use stratified random sampling to create training (70%, 1404 samples), validation (15%, 290 samples), and test (15%, 317 samples) splits. Stratification is performed on the anonymized user number (i.e., ID) to prevent identity leakage across data splits. The same splits are used for all experiments to enable fair comparison.

## A.8. Metadata

Each entry in the RECRUITVIEW metadata file follows the structure shown below. Note that personally identifiable information (user name) has been anonymized.

```
1   {
2       "id": "0001",
3       "video_id": "vid_0001",
4       "video_filename": "vid_0001.mp4",
5       "duration": "long",
6       "question_id": "1",
7       "question": "Introduce yourself",
8       "video_quality": "High",
9       "user_no": "147",
10      "Openness (O)": -0.653,
11      "Conscientiousness (C)": -0.049,
12      "Extraversion (E)": -0.691,
13      "Agreeableness (A)": -0.293,
14      "Neuroticism (N)": 0.190,
15      "overall_personality": -0.029,
16      "interview_score": -0.923,
17      "answer_score": -0.803,
18      "speaking_skills": -0.769,
19      "confidence_score": -0.362,
20      "facial_expression": -0.817,
21      "overall_performance": -0.456,
22      "transcript": "[00:01 - 00:11] Hello everyone, this is ..."
23  }
```

The twelve continuous scores are normalized and represent relative performance across the dataset, derived from the nuclear-norm regularized MNL model described in Section 3.3.

## B. Ethics

### B.1. Participant Protection and Data Collection

All data collection procedures for the RECRUITVIEW dataset were conducted under institutional ethical approval and followed human research standards consistent with the Declaration of Helsinki. Participants were fully briefed about the study's purpose and provided written informed consent prior to participation. The consent form explicitly covered the recording of interview videos, data usage for academic research, and the voluntary nature of participation. Participants were informed of their right to withdraw their data at any point before public release of the dataset used in this study, without consequence. No personally identifiable information (PII) was stored alongside the recordings. All metadata were anonymized, and the participant entries are linked only to an anonymized user number.

Participants were recruited through voluntary university outreach programs and online calls for participation. The pool consisted primarily of adult volunteers, with no inclusion of vulnerable populations.

### B.2. Data Annotation and Labeling

Annotations were performed by clinical psychologists familiar with behavioral and personality assessment protocols. A pairwise comparison framework was adopted instead of absolute rating to reduce inter-rater calibration bias and to ensure consistency across annotators. Comparative judgments were aggregated using a nuclear-norm regularized multinomial logit model to derive continuous, psychometrically consistent target scores. Annotators were compensated fairly for their professional effort. All annotator identities remain confidential.

### B.3. Dataset and Model: Bias, Misuse, and Fairness

The dataset's participant pool exhibits diversity in gender, accent, and educational background, but full demographic uniformity is not available. As such, models trained on this dataset may not generalize equally across other population subgroups. We explicitly acknowledge this limitation and encourage future fairness audits. The RECRUITVIEW dataset and the associated CRMF model are intended solely for research on multimodal behavioral and personality assessment. They are not validated for real-world deployment, hiring processes, or psychological diagnostics. Any attempt to use this work for employment screening, psychological profiling, or commercial analytics constitutes misuse. Although care was taken to minimize annotation bias and maintain fairness, all data-driven systems may still inherit spurious correlations; future work will include comprehensive fairness and subgroup analyses.

### B.4. Responsible Research Practices

We emphasize transparency regarding dataset scope and limitations, including moderate dataset size and short interview durations. The dataset and code are released for non-commercial academic research to enable independent verification, benchmarking, and fairness assessment by the broader community.

Access to the RECRUITVIEW dataset is managed through a secure request portal. Applicants must submit an access request and sign a data usage agreement confirming (1) exclusive non-commercial academic use, (2) adherence to participant anonymity, and (3) compliance with the ethical guidelines outlined in this paper. The agreement prohibits using the dataset or models for employment decisions, identity profiling, or any commercial product development. Each access request is manually reviewed, and credentials are issued only upon approval and formal consent acknowledgment. Access logs are maintained, and the authors reserve the right to revoke access in cases of misuse.

The dataset is released under the license CC BY-NC 4.0 restricting commercial use. All data management and access mechanisms comply with

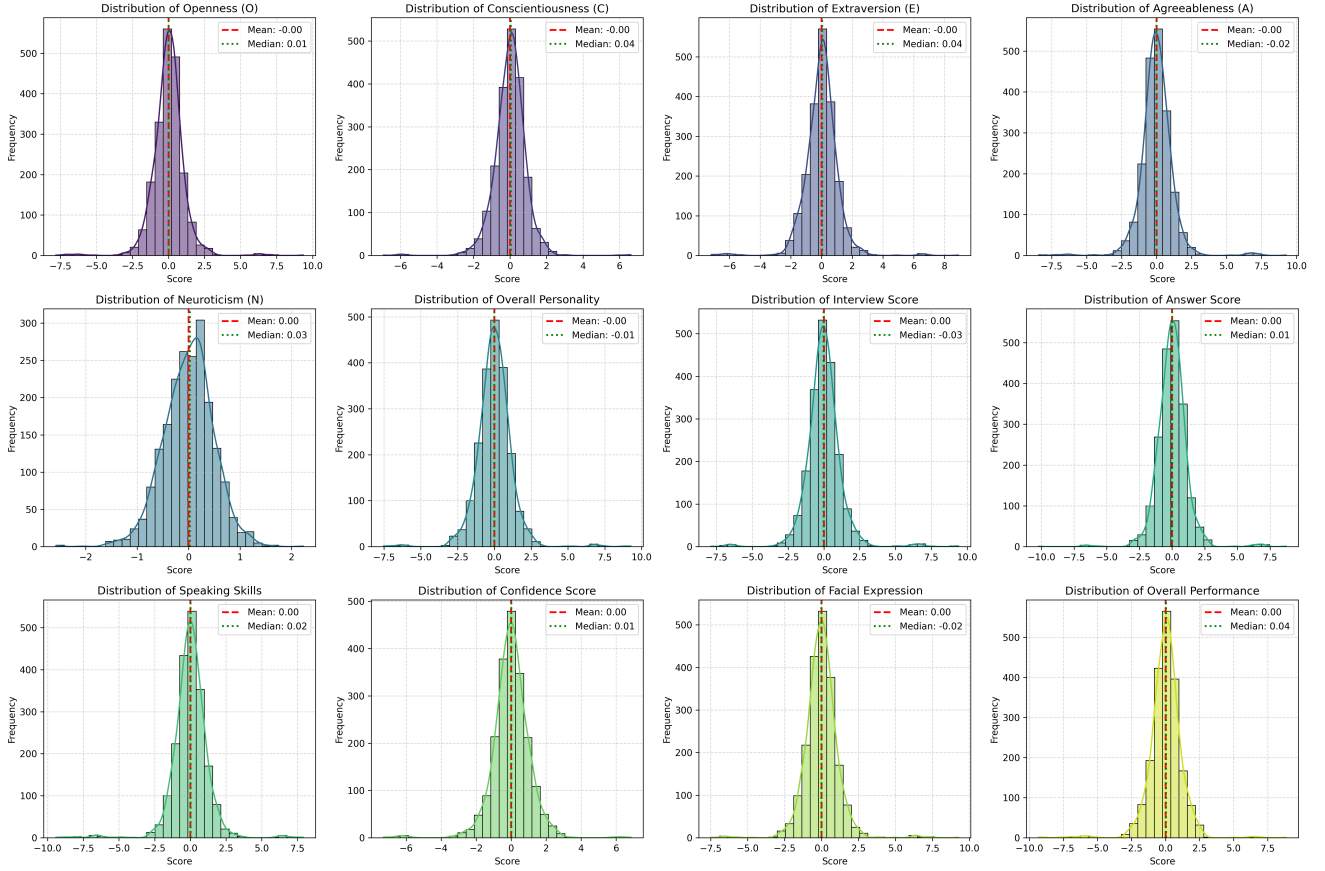Distribution of Target Metrics in the RECRUITVIEW Dataset

Figure 10. Distribution histograms for all 12 target metrics in RECRUITVIEW. Each metric is normalized with a mean near zero. The plots show varying degrees of skewness and heavy tails (leptokurtosis), motivating the use of robust loss functions and rank-based evaluation.

| Statistic | O | C | E | A | N | Pers. | Int. | Ans. | Spk. | Conf. | Fac. | Perf. |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| count | 2011.000 | 2011.000 | 2011.000 | 2011.000 | 2011.000 | 2011.000 | 2011.000 | 2011.000 | 2011.000 | 2011.000 | 2011.000 | 2011.000 |
| mean | -0.000 | -0.000 | -0.000 | -0.000 | -0.000 | -0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 |
| std | 1.127 | 0.877 | 1.098 | 1.200 | 0.487 | 1.203 | 1.230 | 1.176 | 1.278 | 1.080 | 1.151 | 1.201 |
| min | -7.080 | -6.929 | -7.184 | -8.412 | -2.577 | -7.526 | -7.860 | -10.199 | -9.406 | -7.286 | -7.461 | -9.311 |
| 25% | -0.540 | -0.448 | -0.546 | -0.566 | -0.306 | -0.611 | -0.596 | -0.611 | -0.558 | -0.535 | -0.584 | -0.575 |
| 50% | 0.009 | 0.043 | 0.044 | -0.016 | 0.026 | -0.011 | -0.027 | 0.006 | 0.020 | 0.011 | -0.018 | 0.043 |
| 75% | 0.522 | 0.485 | 0.544 | 0.548 | 0.303 | 0.617 | 0.545 | 0.579 | 0.616 | 0.553 | 0.549 | 0.601 |
| max | 9.339 | 6.617 | 8.910 | 9.242 | 2.226 | 9.270 | 9.386 | 8.722 | 7.901 | 6.839 | 9.256 | 8.049 |
| mode | -7.080 | -6.929 | -7.184 | -8.412 | -2.577 | -7.526 | -7.060 | -10.199 | -9.406 | -7.286 | -7.461 | -9.311 |
| skew | 0.027 | -0.570 | -0.218 | 0.398 | -0.245 | 0.462 | 0.660 | 0.353 | -0.855 | -0.640 | 0.598 | -0.748 |
| kurtosis | 12.408 | 10.763 | 11.960 | 13.445 | 1.140 | 10.499 | 11.604 | 12.033 | 12.727 | 8.821 | 11.179 | 11.742 |

Table 7. Comprehensive statistical summary of all 12 target dimensions in RECRUITVIEW. The table shows distribution statistics including central tendency, dispersion, range, and shape measures for the Big Five personality traits (O=Openness, C=Conscientiousness, E=Extraversion, A=Agreeableness, N=Neuroticism), Overall Personality (Pers.), and six interview performance metrics (Int.=Interview Score, Ans.=Answer Score, Spk.=Speaking Skills, Conf.=Confidence Score, Fac.=Facial Expression, Perf.=Overall Performance). Near-zero means confirm proper normalization, while the skewness and kurtosis values indicate the presence of outliers and heavy tails in some dimensions.

institutional data-protection policies and relevant data privacy regulations (e.g., GDPR).

## B.5. Risk and Mitigation Statement

While the RECRUITVIEW dataset contributes valuable insights into multimodal human behavior, we acknowledge potential societal risks. Automated evaluation models trained on human behavioral data could be misinterpreted as objective assessment tools. To mitigate such risks, we provide explicit usage guidelines, controlled data access, and emphasize the need for human oversight in any interpretive use. Continuous monitoring of dataset access and transparency in documentation are maintained to minimize misuse and promote ethical research practice.

## C. Detailed Results

### C.1. Complete Per-Trait and Per-Dimension Results

**Personality Trait Analysis:** Openness shows the strongest CRMF performance ($\rho = 0.6384$ for VMAE+w2v2), representing a 13.4% improvement over the best baseline. This trait measures intellectual curiosity, creativity, and preference for novelty, which likely manifest through diverse behavioral cues across modalities. Conscientiousness, Extraversion, and Agreeableness exhibit moderate but consistent improvements (8-13% gains). Neuroticism presents the most challenging prediction task, with all models achieving lower correlations, though CRMF still improves upon baselines by 24.4%.

**Performance Dimension Analysis:** Interview Score and Answer Score show the strongest absolute performance, with 9-12% improvements over baselines. These metrics directly assess overall interview quality and content quality, which benefit from comprehensive multimodal analysis. Speaking Skills and Confidence Score achieve moderate but consistent improvements (10-16% gains). Overall Performance benefits most from CRMF ($\rho = 0.6521$ for VMAE+HuB), with 9.1% improvement over the strongest baseline.

### C.2. Complete Ablation Study Results

**Fusion Strategy:** Using only simple concatenation or weighted averaging causes severe degradation (21.8% and 17.9% drops), demonstrating that naive fusion strategies fail to capture complex multimodal relationships.

**Pre-Fusion Module:** Removing pre-fusion cross-modal attention or replacing attention pooling with mean pooling results in moderate performance loss (6.7% and 5.6%), confirming that early cross-modal integration provides valuable information flow.

**Geometry Ablations:** Using only a single geometric space consistently underperforms the full model. No single geometry matches the full model, confirming that different geometric spaces capture complementary information. Combining two geometries improves upon single-geometry variants, but all still underperform the full model (3.4% and 3.2% drops for the best two-geometry variants).

**Routing Mechanism:** Hard routing performs comparably to soft routing with only 2.4% degradation. However, removing the router entirely and using uniform weights causes substantial degradation (8.3% drop).

**Prediction Head:** Replacing attention pooling with mean pooling causes severe degradation (14.1% drop). Using a simple linear head instead of the parameter-efficient architecture also substantially degrades performance (11.8% drop). Participants were informed of their right to withdraw

**Loss Function:** Using fixed loss weights reduces performance by 9.2%. Using MSE only without correlation and covariance losses causes 10.1% degradation.

**Architectural Simplifications:** Removing manifold projections causes 4.0% degradation. Removing expert processing entirely leads to 5.7% drop.

**Single Modality:** Video provides the strongest unimodal signal ($\rho = 0.4516$), followed by text ($\rho = 0.4247$) and audio ($\rho = 0.3792$).

## D. Detailed CRMF Architecture

### D.1. Multimodal Encoding

#### D.1.1. Text Encoding Details

We employ DeBERTa-v3-base [36] as our text encoder, which has shown strong performance on natural language understanding tasks. Given tokenized input $\mathbf{T}_{tok} \in \mathbb{Z}^{N_t}$ with attention mask $\mathbf{M}_t \in \{0, 1\}^{N_t}$, the encoder produces contextualized token representations:

$$\mathbf{H}_t = \text{DeBERTa}(\mathbf{T}_{tok}, \mathbf{M}_t) \in \mathbb{R}^{N_t \times d} \qquad (8)$$

where $d = 768$ is the hidden dimension. We fine-tune the last few layers of DeBERTa while keeping earlier layers frozen to balance expressiveness and parameter efficiency. A learned linear projection maps the output to our unified representation space of dimension $d_{model} = 768$.

#### D.1.2. Audio Encoding Details

For audio processing, we explore two self-supervised speech representations: Wav2Vec2 [37] and HuBERT [38]. Given raw audio waveform $\mathbf{A} \in \mathbb{R}^L$ sampled at 16kHz, the encoder produces frame-level representations:

$$\mathbf{H}_a = \text{AudioEncoder}(\mathbf{A}) \in \mathbb{R}^{N_a \times d} \qquad (9)$$

Both Wav2Vec2 and HuBERT learn rich acoustic representations through contrastive predictive coding and masked prediction objectives, respectively. We fine-tune the last few transformer layers while keeping the convolutional feature extractor fixed. The output is projected to $d_{model}$ dimensions.

| Model | Openness (O) | | | Conscientiousness (C) | | | Extraversion (E) | | | Agreeableness (A) | | | Neuroticism (N) | | | Overall Personality | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | $\rho$ | $\tau$-b | C-idx | $\rho$ | $\tau$-b | C-idx | $\rho$ | $\tau$-b | C-idx | $\rho$ | $\tau$-b | C-idx | $\rho$ | $\tau$-b | C-idx | $\rho$ | $\tau$-b | C-idx |
| MiniCPM-o 2.6 (8B) | 0.5629 | 0.3893 | 0.6955 | 0.5014 | 0.3475 | 0.6746 | 0.4996 | 0.3428 | 0.6723 | 0.5484 | 0.3819 | 0.6917 | 0.2378 | 0.1611 | 0.5814 | 0.5613 | 0.3912 | 0.6964 |
| VideoLLaMA2.1-AV (7B) | 0.5474 | 0.3799 | 0.6928 | 0.5077 | 0.3567 | 0.6813 | 0.5021 | 0.3457 | 0.6757 | 0.5397 | 0.3795 | 0.6926 | 0.2252 | 0.1522 | 0.5790 | 0.5309 | 0.3677 | 0.6867 |
| Qwen2.5-Omni (7B) | 0.5354 | 0.3679 | 0.6808 | 0.4957 | 0.3447 | 0.6693 | 0.4901 | 0.3337 | 0.6637 | 0.5277 | 0.3675 | 0.6806 | 0.2132 | 0.1402 | 0.5670 | 0.5189 | 0.3557 | 0.6747 |
| CRMF (VMAE + w2v2) | **0.6384** | **0.4524** | **0.7410** | 0.5572 | **0.4019** | 0.7157 | 0.5681 | 0.4057 | 0.7176 | **0.5927** | **0.4271** | **0.7283** | 0.2603 | 0.1852 | 0.6075 | **0.6098** | **0.4387** | **0.7341** |
| CRMF (VMAE + HuB) | 0.6120 | 0.4312 | 0.7304 | 0.5459 | 0.3889 | 0.7093 | 0.5637 | 0.4019 | 0.7158 | 0.5930 | 0.4231 | 0.7263 | **0.2958** | **0.2101** | **0.6199** | 0.6079 | 0.4371 | 0.7333 |
| CRMF (TimeS + w2v2) | 0.6317 | 0.4480 | 0.7353 | **0.5581** | 0.4000 | **0.7113** | 0.5613 | 0.3954 | 0.7091 | 0.5827 | 0.4166 | 0.7196 | 0.1677 | 0.1211 | 0.5720 | 0.6026 | 0.4314 | 0.7270 |
| CRMF (TimeS + HuB) | 0.6171 | 0.4322 | 0.7299 | 0.5373 | 0.3822 | 0.7049 | **0.5809** | **0.4116** | **0.7196** | 0.5789 | 0.4132 | 0.7204 | 0.2900 | 0.2082 | 0.6180 | 0.6081 | 0.4353 | 0.7315 |

Table 8. Per-trait personality assessment results. CRMF substantially outperforms baselines across all Big Five dimensions and overall personality score. Neuroticism shows the most challenging prediction pattern, consistent with its complex behavioral manifestations. Best results per trait are bolded.

| Model | Interview Score | | | Answer Score | | | Speaking Skills | | | Confidence Score | | | Facial Expression | | | Overall Performance | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | $\rho$ | $\tau$-b | C-idx | $\rho$ | $\tau$-b | C-idx | $\rho$ | $\tau$-b | C-idx | $\rho$ | $\tau$-b | C-idx | $\rho$ | $\tau$-b | C-idx | $\rho$ | $\tau$-b | C-idx |
| MiniCPM-o 2.6 (8B) | 0.5682 | 0.3954 | 0.6985 | 0.5310 | 0.3722 | 0.6869 | 0.5247 | 0.3635 | 0.6826 | 0.5261 | 0.3626 | 0.6821 | 0.4634 | 0.3212 | 0.6614 | 0.5978 | 0.4200 | 0.7108 |
| VideoLLaMA2.1-AV (7B) | 0.5627 | 0.3969 | 0.7013 | 0.5218 | 0.3743 | 0.6900 | 0.5216 | 0.3666 | 0.6862 | 0.5052 | 0.3490 | 0.6774 | 0.4602 | 0.3232 | 0.6645 | 0.5777 | 0.4054 | 0.7056 |
| Qwen2.5-Omni (7B) | 0.5507 | 0.3849 | 0.6893 | 0.5098 | 0.3623 | 0.6780 | 0.5096 | 0.3546 | 0.6742 | 0.4932 | 0.3370 | 0.6654 | 0.4482 | 0.3112 | 0.6525 | 0.5657 | 0.3934 | 0.6936 |
| CRMF (VMAE + w2v2) | 0.6246 | 0.4488 | 0.7392 | 0.5953 | 0.4298 | 0.7297 | **0.5947** | **0.4242** | **0.7269** | 0.5898 | 0.4196 | 0.7246 | 0.5355 | 0.3800 | 0.7049 | 0.6519 | 0.4697 | 0.7496 |
| CRMF (VMAE + HuB) | 0.6180 | 0.4399 | 0.7347 | 0.5919 | 0.4264 | 0.7279 | 0.5804 | 0.4140 | 0.7217 | **0.5950** | **0.4204** | **0.7249** | 0.5179 | 0.3635 | 0.6966 | **0.6521** | 0.4674 | 0.7484 |
| CRMF (TimeS + w2v2) | **0.6299** | **0.4496** | 0.7361 | **0.5968** | **0.4316** | **0.7271** | 0.5894 | 0.4209 | 0.7218 | 0.5903 | 0.4164 | 0.7195 | **0.5387** | **0.3815** | **0.7021** | 0.6477 | 0.4633 | 0.7429 |
| CRMF (TimeS + HuB) | 0.6249 | 0.4438 | **0.7357** | 0.5925 | 0.4227 | 0.7252 | 0.5873 | 0.4180 | 0.7228 | 0.6112 | 0.4309 | 0.7293 | 0.5173 | 0.3623 | 0.6950 | 0.6507 | **0.4639** | **0.7457** |

Table 9. Per-dimension performance assessment results. CRMF shows substantial improvements across all performance metrics, particularly for interview evaluation and overall performance scoring. Facial expression remains challenging but shows consistent gains. Best results per dimension are bolded.

| Component | Variant | Spearman $\rho$ | Kendall $\tau$-b | C-index | Pearson $r$ | MSE |
|---|---|---|---|---|---|---|
| *Full CRMF Model* | | *0.5682* | *0.4069* | *0.7183* | *0.5475* | *0.6864* |
| Fusion Only | Simple Concatenation | 0.4441 | 0.2903 | 0.6151 | 0.4221 | 0.9306 |
| | Weighted Average | 0.4664 | 0.3078 | 0.6239 | 0.4321 | 0.8265 |
| Pre-Fusion | Mean Pooling | 0.5365 | 0.3802 | 0.7026 | 0.5137 | 0.7657 |
| | No Pre-Fusion | 0.5304 | 0.3745 | 0.6972 | 0.5096 | 0.7185 |
| Single Geometry | Hyperbolic Only | 0.5080 | 0.3611 | 0.6980 | 0.4645 | 0.7745 |
| | Spherical Only | 0.5338 | 0.3808 | 0.7054 | 0.4083 | 0.8170 |
| | Euclidean Only | 0.5284 | 0.3753 | 0.7001 | 0.4922 | 0.7526 |
| Two Geometries | Hyperbolic + Spherical | 0.5489 | 0.3916 | 0.7108 | 0.5170 | 0.7403 |
| | Hyperbolic + Euclidean | 0.5161 | 0.3752 | 0.7076 | 0.4583 | 0.8935 |
| | Spherical + Euclidean | 0.5502 | 0.3892 | 0.6993 | 0.5239 | 0.7108 |
| Routing | Hard Routing | 0.5548 | 0.3955 | 0.7127 | 0.5403 | 0.6998 |
| | Uniform Weights (No Router) | 0.5209 | 0.3688 | 0.6994 | 0.4989 | 0.9859 |
| Prediction Head | Mean Pooling (No Attention) | 0.4878 | 0.3527 | 0.7063 | 0.4246 | 0.8139 |
| | Simple Linear Head | 0.5013 | 0.3569 | 0.6984 | 0.4484 | 0.8050 |
| Loss Function | Fixed Loss Weights | 0.5160 | 0.3708 | 0.7104 | 0.4890 | 0.7509 |
| | MSE Only | 0.5110 | 0.3675 | 0.7087 | 0.4885 | 0.7570 |
| Architecture | Linear Projection (No Manifolds) | 0.5457 | 0.3869 | 0.7059 | 0.5294 | 0.7440 |
| | No Expert Processing | 0.5360 | 0.3828 | 0.7014 | 0.4877 | 0.7713 |
| Single Modality | Video Only (VMAE) | 0.4516 | 0.2974 | 0.6521 | 0.4138 | 0.8847 |
| | Audio Only (Wav2Vec2) | 0.3792 | 0.2461 | 0.6245 | 0.3482 | 1.0516 |
| | Text Only (DeBERTa) | 0.4247 | 0.2806 | 0.6429 | 0.3895 | 0.9324 |

Table 10. Complete systematic ablation study results. All experiments use VideoMAE+Wav2Vec2 encoders.

### D.1.3. Video Encoding with Temporal Modeling

For visual processing, we investigate two video understanding architectures: VideoMAE [39] and TimeSformer [40]. Given input video with variable frame count, we first apply 3D convolutional interpolation to adapt the temporal dimension to the encoder's expected frame count (16 for VideoMAE, 8 for TimeSformer). For an input $\mathbf{V} \in \mathbb{R}^{T \times 3 \times 224 \times 224}$, this yields $\mathbf{V}' \in \mathbb{R}^{T' \times 3 \times 224 \times 224}$.

The encoder extracts patch-level features, which we re-shape into temporal-spatial structure. We apply spatial average pooling to obtain temporal features $\mathbf{F}_v \in \mathbb{R}^{T' \times d}$.

To capture rich temporal dynamics, we apply a multi-stage temporal modeling pipeline:

$$\mathbf{F}_{lstm} = \text{BiLSTM}(\mathbf{F}_v) \in \mathbb{R}^{T' \times d} \quad (10)$$

$$\mathbf{F}_{attn} = \text{MultiHeadAttn}(\mathbf{F}_{lstm}\mathbf{F}_{lstm}, \mathbf{F}_{lstm}) \in \mathbb{R}^{T' \times d} \quad (11)$$

$$\mathbf{F}_{conv} = \text{Conv1D}(\mathbf{F}_{attn}) \in \mathbb{R}^{T' \times d} \quad (12)$$

$$\mathbf{H}_v = \text{Proj}(\mathbf{F}_{lstm} + \mathbf{F}_{attn} + \mathbf{F}_{conv}) \in \mathbb{R}^{T' \times d_{model}} \quad (13)$$

where BiLSTM captures sequential dependencies, multihead attention models long-range interactions, and depthwise convolution captures local temporal patterns. The multi-scale fusion combines all three views, and a final projection maps to $d_{model}$ dimensions. This produces a temporal sequence $\mathbf{H}_v \in \mathbb{R}^{8 \times 768}$ preserving fine-grained temporal information for subsequent pre-fusion processing.

VideoMAE employs masked autoencoding with high masking ratios for efficient self-supervised learning, while TimeSformer uses divided space-time attention. We fine-tune the last few transformer blocks of each encoder while keeping earlier layers frozen for parameter efficiency.

### D.2. Pre-Fusion Module

The pre-fusion module performs early integration of multimodal features through cross-modal attention. We concate-

nate encoded features from all modalities and add learnable modality embeddings to distinguish information sources:

$$\mathbf{H}_{cat} = [\mathbf{H}_t; \mathbf{H}_a; \mathbf{H}_v] + \mathbf{E}_{mod} \tag{14}$$

where $\mathbf{E}_{mod} \in \mathbb{R}^{3 \times d}$ contains unique embeddings for text, audio, and video. A multi-layer transformer encoder with $L_{pre}$ layers processes the concatenated sequence:

$$\mathbf{H}_{fused} = \text{Transformer}_{pre}(\mathbf{H}_{cat}) \tag{15}$$

enabling rich cross-modal interactions through self-attention.

To obtain a fixed-dimensional clip-level representation, we employ learned attention pooling rather than simple mean pooling:

$$\alpha_i = \frac{\exp(\mathbf{w}^\top \mathbf{h}_i)}{\sum_j \exp(\mathbf{w}^\top \mathbf{h}_j)}, \quad \mathbf{z}_{pre} = \sum_i \alpha_i \mathbf{h}_i \tag{16}$$

where $\mathbf{w} \in \mathbb{R}^d$ is a learnable attention vector and $\mathbf{h}_i$ denotes the $i$-th token in $\mathbf{H}_{fused}$. This pooling mechanism learns to emphasize tokens most relevant for behavioral assessment, producing $\mathbf{z}_{pre} \in \mathbb{R}^d$.

### D.3. Geometric Expert Architectures

#### D.3.1. Hyperbolic Expert Mathematical Details

The hyperbolic expert performs operations in the gyrovector space framework [42], using Möbius transformations that preserve hyperbolic distances. For a $L_{exp}$-layer network:

$$\mathbf{x}_h^{(\ell+1)} = \sigma_h \left( \mathbf{W}_h^{(\ell)} \otimes_c \mathbf{x}_h^{(\ell)} \oplus_c \mathbf{b}_h^{(\ell)} \right) \tag{17}$$

where $\otimes_c$ denotes Möbius matrix-vector multiplication, $\oplus_c$ is Möbius addition, and $\sigma_h$ is a Möbius pointwise nonlinearity. Specifically, Möbius addition is defined as:

$$\mathbf{x} \oplus_c \mathbf{y} = \frac{(1 + 2c\langle\mathbf{x},\mathbf{y}\rangle + c\|\mathbf{y}\|^2)\mathbf{x} + (1 - c\|\mathbf{x}\|^2)\mathbf{y}}{1 + 2c\langle\mathbf{x},\mathbf{y}\rangle + c^2\|\mathbf{x}\|^2\|\mathbf{y}\|^2} \tag{18}$$

The Möbius pointwise nonlinearity applies activation functions in tangent space: $\sigma_h(\mathbf{x}) = \exp_\mathbf{x}^c(\sigma(\log_\mathbf{x}^c(\mathbf{x})))$, where $\log_\mathbf{x}^c$ and $\exp_\mathbf{x}^c$ are the logarithmic and exponential maps at $\mathbf{x}$.

After processing, we apply residual connections using Möbius addition: $\mathbf{x}_h^{out} = \mathbf{x}_h^{(L)} \oplus_c \mathbf{x}_h^{(0)}$. All operations preserve the hyperbolic geometry, ensuring outputs remain in the Poincaré ball.

#### D.3.2. Spherical Expert Mathematical Details

Operations on the sphere are performed in tangent space via exponential and logarithmic maps. Given base point $\mathbf{p}$ (we

use the north pole), the logarithmic map projects $\mathbf{x}_s$ to the tangent space $T_\mathbf{p}\mathbb{S}^{d-1}$:

$$\log_\mathbf{p}(\mathbf{x}_s) = \frac{\arccos(\langle\mathbf{p},\mathbf{x}_s\rangle)}{\sqrt{1 - \langle\mathbf{p},\mathbf{x}_s\rangle^2}}(\mathbf{x}_s - \langle\mathbf{p},\mathbf{x}_s\rangle\mathbf{p}) \tag{19}$$

In tangent space, standard linear transformations and activations apply:

$$\mathbf{v}^{(\ell+1)} = \sigma(\mathbf{W}_s^{(\ell)}\mathbf{v}^{(\ell)} + \mathbf{b}_s^{(\ell)}) \tag{20}$$

where $\mathbf{v}^{(\ell)} \in T_\mathbf{p}\mathbb{S}^{d-1}$. The final tangent vector is mapped back to the sphere via exponential map:

$$\exp_\mathbf{p}(\mathbf{v}) = \cos(\|\mathbf{v}\|)\mathbf{p} + \sin(\|\mathbf{v}\|)\frac{\mathbf{v}}{\|\mathbf{v}\|} \tag{21}$$

Residual connections in tangent space combine the input and output: $\mathbf{v}^{out} = \mathbf{v}^{(L)} + \mathbf{v}^{(0)}$, followed by exponential map back to $\mathbb{S}^{d-1}$.

#### D.3.3. Euclidean Expert

The Euclidean expert uses standard feed-forward layers with residual connections:

$$\mathbf{x}_e^{(\ell+1)} = \text{ReLU}(\mathbf{W}_e^{(\ell)}\mathbf{x}_e^{(\ell)} + \mathbf{b}_e^{(\ell)}), \quad \mathbf{x}_e^{out} = \mathbf{x}_e^{(L)} + \mathbf{x}_e^{(0)} \tag{22}$$

Each expert has $L_{exp}$ layers with dropout rate $p$ for regularization.

### D.4. Geometry-Aware Attention Details

To further refine expert outputs, we apply intra-manifold attention that respects geometric structure. For each geometry, we compute attention in its respective tangent space.

#### D.4.1. Hyperbolic Intra-Manifold Attention

Given hyperbolic representations $\mathbf{x}_h^{out}$, we map to tangent space at the origin:

$$\mathbf{v}_h = \log_\mathbf{0}^c(\mathbf{x}_h^{out}) = \frac{\text{arctanh}(\sqrt{c}\|\mathbf{x}_h^{out}\|)}{\sqrt{c}\|\mathbf{x}_h^{out}\|}\mathbf{x}_h^{out} \tag{23}$$

Multi-head self-attention is applied in tangent space (which is Euclidean):

$$\mathbf{v}_h^{att} = \text{MultiHead}(\mathbf{Q}_h, \mathbf{K}_h, \mathbf{V}_h) \tag{24}$$

where $\mathbf{Q}_h = \mathbf{W}_Q\mathbf{v}_h$, $\mathbf{K}_h = \mathbf{W}_K\mathbf{v}_h$, $\mathbf{V}_h = \mathbf{W}_V\mathbf{v}_h$. The attended representation is mapped back:

$$\mathbf{x}_h^{att} = \exp_\mathbf{0}^c(\mathbf{v}_h^{att}) \tag{25}$$

#### D.4.2. Spherical and Euclidean Attention

Similar procedures apply for spherical geometry using $\log_\mathbf{p}$ and $\exp_\mathbf{p}$. For Euclidean space, attention is applied directly without manifold conversions. All attention modules use multiple heads with temperature scaling $\tau$ to sharpen attention distributions.

## D.5. Routing Mechanism Details

### D.5.1. Routing Regularization

To encourage diverse geometry utilization, we apply entropy regularization on routing weights:

$$\mathcal{L}_{entropy} = -\lambda_{ent}H(\mathbf{r}) = -\lambda_{ent}\sum_{i=1}^{K} r_i \log r_i \qquad (26)$$

A negative value encourages high entropy (uniform distribution), promoting complementary geometric views rather than specialization. We also apply load balancing regularization:

$$\mathcal{L}_{balance} = \lambda_{bal}\text{Var}(\mathbb{E}_{batch}[\mathbf{r}]) \qquad (27)$$

ensuring all experts are utilized across the dataset.

## D.6. Geometric Fusion Theoretical Justification

The tangent space fusion strategy is equivalent to first-order Fréchet mean approximation on the product manifold $\mathcal{M} = \mathbb{B}_c^{d_e} \times \mathbb{S}^{d_e-1} \times \mathbb{R}^{d_e}$. The Fréchet mean minimizes:

$$\mu^* = \arg\min_{\mathbf{x}} \sum_i w_i d_{\mathcal{M}_i}^2(\mathbf{x}_i, \mathbf{x}) \qquad (28)$$

where $d_{\mathcal{M}_i}$ is the distance on manifold $\mathcal{M}_i$. The first-order approximation linearizes the problem in tangent space, yielding the weighted combination. This approach avoids expensive iterative optimization while providing a theoretically grounded fusion mechanism.

## D.7. Multi-Task Prediction Head Details

### D.7.1. Shared Representation Learning

A shared MLP processes the fused features:

$$\mathbf{h}^{(1)} = \text{GELU}(\text{LayerNorm}(\mathbf{W}_{sh}^{(1)}\mathbf{z}_{refined} + \mathbf{b}_{sh}^{(1)})) \quad (29)$$

$$\mathbf{h}^{(2)} = \text{GELU}(\text{LayerNorm}(\mathbf{W}_{sh}^{(2)}\mathbf{h}^{(1)} + \mathbf{b}_{sh}^{(2)})) \qquad (30)$$

producing a shared representation $\mathbf{h}^{(2)} \in \mathbb{R}^{512}$ that captures common structure across all targets.

### D.7.2. Task-Specific Adaptation

Each of the $K = 12$ targets has a lightweight adaptation module:

$$\hat{y}_k = \mathbf{w}_k^{(2)\top}\text{GELU}(\mathbf{W}_k^{(1)}\mathbf{h}^{(2)} + \mathbf{b}_k^{(1)}) + b_k \qquad (31)$$

where $\mathbf{W}_k^{(1)} \in \mathbb{R}^{64\times512}$ and $\mathbf{w}_k^{(2)} \in \mathbb{R}^{64}$ are task-specific parameters. This design dramatically reduces parameters compared to full per-task networks while maintaining expressive capacity.

## D.8. Training Objective Details

### D.8.1. Multi-Component Loss

Our training objective combines multiple loss components through adaptive balancing:

$$\mathcal{L}_{total} = \sum_{i=1}^{N} \beta_i \mathcal{L}_i \qquad (32)$$

where $\mathcal{L}_i$ are individual loss components and $\beta_i$ are adaptive weights learned during training.

**Regression Loss:** We use Huber loss for robustness to outliers:

$$\mathcal{L}_{reg} = \frac{1}{K}\sum_{k=1}^{K} \begin{cases} \frac{1}{2}(y_k - \hat{y}_k)^2 & \text{if } |y_k - \hat{y}_k| \leq \delta \\ \delta(|y_k - \hat{y}_k| - \frac{1}{2}\delta) & \text{otherwise} \end{cases}$$
$$(33)$$

with $\delta = 1.0$. This combines MSE's efficiency for small errors with MAE's robustness for large deviations.

**Correlation Boosting Loss:** To encourage predictions that maintain correlation structure with targets:

$$\mathcal{L}_{corr} = \lambda_{corr}\left(1 - \frac{1}{K}\sum_{k=1}^{K} |\rho(\hat{\mathbf{y}}_k, \mathbf{y}_k)|\right) \qquad (34)$$

where $\rho$ denotes Pearson correlation.

**Covariance Alignment Loss:** To match the covariance structure between predictions and targets:

$$\mathcal{L}_{cov} = \lambda_{cov}\|\text{Cov}(\hat{\mathbf{Y}}) - \text{Cov}(\mathbf{Y})\|_F^2 \qquad (35)$$

where $\text{Cov}(\cdot)$ computes the empirical covariance matrix and $\|\cdot\|_F$ is the Frobenius norm.

**Auxiliary Losses:** Routing regularization losses $\mathcal{L}_{entropy}$ and $\mathcal{L}_{balance}$ are added, along with head regularization encouraging small adapter weights.

## D.9. Adaptive Loss Balancing

Rather than fixed weights, we learn to balance loss components adaptively. Each component has a learnable weight $\alpha_i$ and running exponential moving average (EMA) statistics:

$$\mu_i^{(t)} = \gamma\mu_i^{(t-1)} + (1-\gamma)\mathcal{L}_i^{(t)} \qquad (36)$$

Adaptive weights are computed via inverse variance weighting:

$$\beta_i^{adapt} = \frac{1/(\text{Var}(\mathcal{L}_i) + \epsilon)}{\sum_j 1/(\text{Var}(\mathcal{L}_j) + \epsilon)} \qquad (37)$$

then combined with learned weights: $\beta_i = \omega\text{softmax}(\alpha_i) + (1-\omega)\beta_i^{adapt}$. This balancing prevents any single loss from dominating training.