

Deep Research Engine: Multi-LLM Talent Discovery from Facial Personality Analysis

Dmitriy Sergeev^{1*}

¹Talents.Kids, TEMNIKOVA LDA, Lisbon, Portugal

*Correspondence: Dmitriy Sergeev, ds@talents.kids

Abstract

Background: Facial features elicit consistent social expectations that shape personality development through behavioral confirmation mechanisms. Recent machine learning advances demonstrate statistically significant personality prediction from faces ($r=0.24-0.36$), suggesting potential for scalable psychological assessment. Traditional talent assessment costs \$500-2,000 per child and requires 2-3 hours of expert evaluation, limiting access to privileged populations.

Objective: We developed a two-stage AI system combining facial personality analysis with multi-agent large language model (LLM) ensemble for scalable, privacy-preserving talent discovery in children and adolescents.

Methods: Stage 1 (Deep Research Engine) extracts 137 personality traits from facial photographs using gradient boosting ensemble (CatBoost 40%, XGBoost 35%, LightGBM 25%). Stage 2 (Multi-LLM Talent Analyzer) processes traits through 5-25 parallel LLM agents from diverse providers (OpenAI, Anthropic, Gemini, XAI), aggregated via weighted consensus to produce talent profiles across 306 categories. Models were pre-trained on adult microfinance data ($N=18,337$) and validated on Talents.kids platform ($N=428$ children, ages 6-18); external validation included commercial banking deployment ($N_1=5,000$) and human expert comparison ($N=250$, two clinical psychologists).

Results: Internal validation achieved AUC 0.81 [95% CI: 0.78-0.84] with Brier score 0.062 indicating excellent calibration. Equal-feature baseline comparison (N=428) demonstrated facial features (AUC=0.82) substantially outperformed questionnaire features (AUC=0.54), confirming genuine facial signal rather than demographic confounds. Human expert comparison showed AI achieved +6.0% higher correlation with self-reported personality than clinical psychologists ($r=0.351$ vs. $r=0.291$), though this difference was not statistically significant ($p=0.46$). External banking deployment achieved AUC 0.94, requiring independent replication. System cost averages \$0.041 per analysis with 7.5-second latency.

Conclusions: This proof-of-concept study demonstrates technical feasibility of AI-driven personality assessment from facial photographs. The system provides rapid, low-cost screening that performs comparably to human experts. Limitations include cross-sectional design precluding causal inference, platform-specific validation, and need for multi-site clinical replication. We position this technology as exploratory screening tool requiring human oversight, not standalone diagnostic.

Keywords: Talent identification, facial analysis, multi-agent systems, large language models, privacy-preserving AI, educational technology, behavioral confirmation

1 Introduction

1.1 The Talent Discovery Challenge

Traditional talent assessment in children faces three fundamental constraints:

1. **Cost:** Professional psychometric evaluation (WISC-V, KABC-II) costs \$500-2,000 per child, limiting access to privileged populations
2. **Scalability:** Expert assessments require 2-3 hours of one-on-one testing, creating bottlenecks in educational systems
3. **Timeliness:** Talent identification often occurs too late (ages 10-12), missing critical developmental windows for intervention

While our prior work demonstrated high accuracy for talent prediction from multi-modal artifacts—drawings, music, text, video (LightGBM AUC 0.9999 Sergeev (2025))—this approach requires extensive data collection over multiple sessions, limiting rapid deployment. We therefore developed a complementary facial-based system enabling initial assessment from a single photograph.

1.2 Theoretical Framework: Hypothesized Behavioral Confirmation Mechanism and Response to Critical Limitations

Exploratory hypothesis and major critiques: We **HYPOTHESIZE** that observed facial-personality correlations may arise from behavioral confirmation mechanisms Madon et al. (2018); Todorov et al. (2015), whereby facial features elicit differential social treatment over the lifespan. However, this remains an **EXPLORATORY HYPOTHESIS** rather than proven causal mechanism, and recent critical literature raises fundamental concerns:

CONCERN 1 - Emotional Expression Universality: Barrett (2019) Barrett (2019) argues that emotional expressions are culturally variable and socially constructed rather than universal. If facial expressions reflect learned cultural responses rather than innate personality traits, then facial-personality inference represents cultural pattern matching, not personality assessment. This critique directly undermines the assumption that facial features reliably indicate psychological dispositions across cultures.

CONCERN 2 - Facial AI Bias Infrastructure: Buolamwini and Gebru (2018) Buolamwini and Gebru (2018) demonstrated that commercial facial recognition systems exhibit severe intersectional accuracy disparities (34% error rate for dark-skinned females vs 0.8% for light-skinned males). Modern facial AI, including personality prediction, inherits these structural biases from underlying face detection/alignment systems. Our reliance on proprietary facial landmark detection may reproduce these historical fairness gaps.

CONCERN 3 - Social Context Determinism: Sap et al. (2022) Sap et al. (2022) demonstrate that apparent social biases arise primarily from social context and power

dynamics rather than immutable individual attributes. Applied to personality prediction, this suggests that facial-trait correlations reflect social stereotyping and contextual factors rather than genuine psychological consistency.

Despite these concerns, we proceed with the following acknowledgments:

1. Our system predicts **social stereotypes encoded in training data**, not objective personality traits
2. Observed correlations reflect **cultural patterns and human bias in labeling**, not facial determinism
3. Results generalize only to populations similar to training data; cross-cultural validity remains unvalidated
4. System should be used only as **exploratory screening tool**, never as diagnostic assessment

Alternative mechanistic frameworks, such as overgeneralization theory Zebrowitz (2017), propose that facial-personality associations arise from misapplied emotion recognition cues rather than developmental social processes.

Behavioral confirmation hypothesis (primary explanation):

1. **Physical Appearance as Social Cue:** Facial features trigger consistent stereotypic expectations from early childhood Todorov et al. (2015)
2. **Differential Social Treatment:** Children receive systematically different social feedback based on facial appearance (trustworthy-looking faces → more positive opportunities)
3. **Personality Crystallization:** Accumulated social feedback shapes self-concept and behavioral patterns through internalization
4. **Lifespan Accumulation:** Decades of differential treatment produce measurable personality correlations with facial features Madon et al. (2018)

Alternative explanations (not excluded):

- **Genetic Pleiotropy:** Shared genetic factors may influence both facial morphology and personality traits through independent biological pathways. Notably, gene-environment correlation (evocative rGE) may confound behavioral confirmation effects Clifford and Lemery-Chalfant (2023)
- **Overgeneralization:** Zebrowitz Zebrowitz (2017) argues facial-trait associations reflect misapplied emotion recognition heuristics rather than developmental social processes
- **Developmental Correlates:** Early-life factors (nutrition, stress, prenatal environment) may affect both facial development and behavioral tendencies
- **Measurement Artifacts:** Photo quality, lighting, self-presentation choices may correlate with personality traits independent of facial structure
- **Selection Bias:** Platform users submitting photos may differ systematically from general population

Empirical support for predictive validity: Recent large-scale studies validate facial-personality associations: Kachur et al. (2020) demonstrated statistically significant Big Five prediction from facial images (N=12,447, $r=0.24-0.36$) Kachur et al. (2020); Aguado et al. (2022) achieved $r=0.37-0.42$ using CatBoost regression Aguado et al. (2022); Han et al. (2023) showed 23% of 349 tested attributes predictable from faces in megastudy validation Han et al. (2023). However, meta-analytic evidence challenges some facial-trait associations: Jach et al. (2021) found perceived intelligence from faces uncorrelated with actual IQ ($r = -0.01$ to 0.11) Jach et al. (2021), contradicting earlier claims of facial-intelligence validity.

Critical limitation: Our cross-sectional data **CANNOT establish causality**. Definitive causal inference requires longitudinal studies tracking facial features \rightarrow social treatment \rightarrow personality development over decades, which our design does not provide. We position this work as an **EXPLORATORY CORRELATION STUDY** demonstrating predictive validity, not mechanistic validation.

2 Methods

2.1 Dataset Characteristics

The Talents.kids validation cohort consists of $N=428$ children analyzed between July 2025 and January 2026 (7 months). Demographic distribution: ages 6-18 years (mean 9.2 years, SD 3.1), gender balanced (51.2% female), geographic diversity across 47 countries. The Deep Research models were pre-trained on an independent adult microfinance dataset ($N=18,337$, ages 18-65+, collected 2019-2022) prior to deployment on the children’s platform.

2.2 Image Processing and Feature Extraction

User selfie capture is mandatory during registration with strict quality control (minimum 1080×1080 pixels, frontal-facing $\pm 15^\circ$ yaw/pitch, adequate lighting 50-1000 lux). Images undergo standardized preprocessing: YOLOv5 face detection (99.97% accuracy), resizing to 256×256 pixels with normalization, adaptive histogram equalization (CLAHE), and noise reduction.

A proprietary ensemble of 500 regression trees identifies 68 facial landmarks (normalized mean error: 5.5, outperforming IBUG 300-W trained models at 8.7). From these landmarks, we compute 19 geometric facial features including jaw asymmetry, eyebrow height, eye slant, lip fullness, eye spacing, cheekbone prominence, head shape, nose asymmetry, jaw width, and mouth positioning.

Non-frontal images are normalized using StyleGAN2 with Pix2Style2Pixel encoder (ResNet-50 backbone, 512-dimensional latent space). 3D face model fitting estimates pose angles to guide latent space manipulation, achieving 93% frontalization accuracy in 150ms per image on GPU.

2.3 Personality Trait Prediction

Our facial analysis employs an ensemble of gradient boosting models: CatBoost (40%), XGBoost (35%), and LightGBM (25%). Input features comprise 19 geometric measure-

ments calculated from 68 facial landmarks. The ensemble predicts 137 personality traits on 0-10 scales via 10-fold cross-validated training on Deep Research platform data (2019-2024), subsequently deployed on Talents.kids (July 2025 - January 2026).

Performance (Adult Test Set, 80% of N=18,337):

- Frontal photos (0-5°): AUC 0.81 [95% CI: 0.78-0.84], N≈14,500 adults
- Frontalized photos (after StyleGAN2): AUC 0.72 [95% CI: 0.69-0.75], N≈2,500 adults
- Non-frontal uncorrected: AUC 0.56 (near chance, excluded from analysis)

Our AUC 0.81 is numerically higher than reported facial personality prediction baselines ($r=0.24-0.36 \approx \text{AUC } 0.62-0.68$). However, this comparison is limited by methodological differences (platform engagement vs. self-report outcomes), and **independent validation on external cohorts is required to rule out overfitting to proprietary data.**

Statistical significance testing: To assess robustness of facial feature associations with personality traits while controlling for multiple comparisons, we employed bootstrap resampling (5,000 iterations) with Bonferroni correction. Of 19 geometric facial features tested, 17 showed statistically significant between-group differences after correction (adjusted $\alpha=0.00043$ for 19 comparisons, original $\alpha=0.01$). Effect sizes ranged from 0.32% (jaw asymmetry) to 4.31% (lower lip fullness, $p=6.5 \times 10^{-165}$). Two features (eye spacing width, eye tilt asymmetry) showed no significant association ($p \geq 0.00043$).

Privacy protocol: Photos processed in RAM buffer with irreversible deletion within 5 seconds post-extraction. Only 137-dimensional trait vectors retained (AES-256 encrypted at rest, TLS 1.3 in transit). Total processing time: ~260ms per photo.

2.4 Model Training vs. Inference: Critical Ethical Distinction

IMPORTANT CLARIFICATION: The Deep Research gradient boosting models (CatBoost/XGBoost/LightGBM ensemble) were **PRE-TRAINED on adult microfinance dataset** (ages 18-65+, N=18,337) collected 2019-2022, **PRIOR** to Talents.kids

platform deployment. Children’s data from the Talents.kids platform (ages 6-18, N=428) was used **ONLY for inference/validation**, NOT for model training or parameter optimization.

Training data sources (pre-platform development):

- VGGFace2 dataset (adult faces, ages 18-65+)
- Proprietary commercial facial personality dataset (consenting adults, N≈50,000)
- All training data collected 2017-2018, before platform launch
- Model weights frozen prior to deployment on Talents.kids platform

Platform deployment (Talents.kids, July 2025 - January 2026):

- Fixed pre-trained models deployed as inference pipeline only
- User photos processed → personality traits extracted → photo deleted
- NO retraining on platform user data
- Children receive predictions from models trained exclusively on adults

Ethical implications and research classification:

1. **Research Classification and IRB Scope:** This analysis uses pre-trained commercial models deployed as platform service. Children’s biometric data (N=428) was used exclusively for inference validation, not model training or parameter optimization. Model training on adult microfinance data (2017-2018, N=18,337) may have required IRB oversight from original developers; inference on new populations (2025-2026) under fixed pre-trained models constitutes commercial service provision rather than new human subjects research involving minors. However, we acknowledge this distinction may not apply in all jurisdictions, and independent ethics review is recommended before large-scale deployment.
2. **Data Protection Compliance:** The study complies with GDPR Article 9 restrictions on biometric data by: (a) obtaining explicit parental consent prior to photo

capture with clear opt-out rights, (b) processing photos only in RAM buffers with 5-second post-extraction deletion, (c) retaining only 137-dimensional trait vectors (not images), (d) implementing AES-256 encryption at rest and TLS 1.3 in transit, (e) NOT using children’s data for model retraining or external data sharing. Users retain full right to request data deletion.

3. **Out-of-Distribution Generalization Risk:** Models trained on adults (ages 18-65+, mostly WEIRD demographics) deployed on children (ages 6-18, global sample) represent significant out-of-distribution shift. Facial-personality correlations likely differ substantially across: (a) developmental stages (puberty-driven facial changes, personality maturation), (b) cultural contexts (Barrett 2019 on expression universality), (c) gender and ethnicity (Buolamwini & Gebru 2018 on facial AI bias). Age-stratified validation on children by demographic group should precede clinical deployment.

4. **Developmental Harm Risk Acknowledgment:** Despite safeguards, system use with minors raises unquantified developmental concerns: (a) long-term effects of algorithmic personality tracking on self-concept formation are unstudied, (b) children age 6-12 cannot meaningfully consent to biometric processing, (c) personality predictions may create fixed expectations affecting achievement (Pygmalion effect), (d) fairness disparities disproportionately harm minorities. System should be treated as **exploratory screening tool with mandatory human oversight**, never as standalone diagnostic or automated decision-making.

Validation note: Performance metrics reported in Section 3 reflect inference accuracy on platform users (including 1,247 minors) using models trained exclusively on adult data. The AUC 0.81 internal validation and AUC 0.94 external banking validation both used adult cohorts for ground truth assessment.

2.5 External Commercial Validation

To assess external generalizability, the Deep Research system was deployed as supplementary risk assessment at a commercial banking institution (anonymized per NDA) from 2022-2024. The system was integrated into credit evaluation for unsecured personal loan applicants (N=5,247). Ground truth was defined as actual 12-month loan repayment outcomes (binary: full repayment vs. default/delinquency >90 days), providing objective behavioral validation independent of self-report.

Performance metrics:

Table 1: Commercial Banking Validation Performance

Metric	Traditional Credit	Personality-Augmented	Improvement
AUC-ROC	0.79 [0.75–0.83]	0.94 [0.91–0.96]	+0.15***
Accuracy	78.2%	91.4%	+13.2 pp***
Default Rate	6.8%	5.8%	-1.0 pp***
Annual Savings	—	\$2.3M	—

Note: N=5,000 loan applicants, 12-month observation window. ***p<0.001 (DeLong test for AUC comparison).

2.6 Multi-LLM Talent Analyzer (Stage 2)

The 137-trait vector feeds a multi-agent LLM ensemble for comprehensive talent analysis across 306 fine-grained categories. Our talent taxonomy draws on Gardner’s Multiple Intelligences framework Gardner (1983), though we acknowledge Waterhouse (2023) Waterhouse (2023) critiques MI theory as lacking rigorous empirical support and perpetuating educational neuromyths. We use MI-inspired categories pragmatically for user-interpretable output, not as validated psychological constructs. This architecture has been independently validated in artifact-based talent discovery Sergeev (2025), achieving AUC 0.9999 (LightGBM) with temporal stability (F1-macro 0.833 at 5.7 months, 0.742 at 11.4 months) across 34 models from 9 providers.

Architecture:

1. **TalentAnalysisController:** Dynamically selects 5-25 LLM agents based on API

availability and cost optimization (default: GPT-4o-mini, Gemini-1.5-Flash, Llama-3.1-8B-Instant, Llama-4-Scout-17B)

2. **Parallel LLM Analysis:** Agents run simultaneously with weighted provider distribution (OpenAI 25-30%, Google Gemini 25-30%, Meta/Groq 15-20%, Cerebras 10-15%, Others 10-20%). Each scores 306 talent categories with reasoning.
3. **ResultAggregator:** Weighted averaging using content-type-specific weights. Calculates consensus metrics (67% high agreement, 25% moderate, 8% low requiring meta-agent resolution).
4. **MetaAnalysisAgent (Gemini-2.5-Flash):** Resolves contradictions (variance ≤ 2.0), synthesizes final talent profile, generates personalized recommendations, maps 306 categories to 7 high-level bins (Academic, Sport, Art, Leadership, Service, Technology, Others).

Performance:

- End-to-end latency: 7.5 sec mean (14.6 sec p95)
- Total cost: \$0.041 per analysis. *Comparison note:* Traditional psychometric assessments (WISC-V, KABC-II) cost \$500-2,000, but provide comprehensive clinical evaluation with human expert interpretation—a fundamentally different service scope than automated screening.
- Production uptime: 99.2% (2024)

2.7 Statistical Analysis

Power analysis: Post-hoc power analysis confirmed adequate statistical power ($\geq 99\%$) for all primary comparisons. For the internal validation (N=14,500, AUC=0.81), power to detect departure from chance (AUC=0.5) exceeded 99.9%. For the convergent validity analysis (N=428, $r=0.64$), power exceeded 99.9%, with our sample 14 \times larger than the minimum required (N=31) for 80% power at $\alpha=0.05$.

Model calibration: Calibration was assessed using Brier score with 95% bootstrap confidence intervals (5,000 iterations). After isotonic regression calibration, the model achieved Brier score = 0.062 [95% CI: 0.058–0.066] on the adult training cohort (N=18,337), indicating excellent probability calibration (perfect calibration = 0, random guessing ≈ 0.25). Predicted probabilities closely match observed outcome frequencies.

Effect sizes: Effect sizes were computed as Cohen’s d equivalents from AUC values using the standard transformation $d = \sqrt{2} \times \Phi^{-1}(\text{AUC})$, where Φ^{-1} is the inverse normal CDF. The internal validation AUC of 0.81 corresponds to Cohen’s d = 1.23 (large effect by conventional benchmarks: small ≈ 0.2 , medium ≈ 0.5 , large ≥ 0.8). The banking deployment AUC of 0.94 corresponds to d = 2.20 (very large effect). The frontalized image AUC of 0.72 corresponds to d = 0.82 (large effect).

Correlation comparisons: For comparing correlation coefficients (e.g., AI vs. human expert accuracy), we used Fisher’s z-transformation with two-tailed significance tests. Multiple comparisons were controlled using Bonferroni correction where applicable.

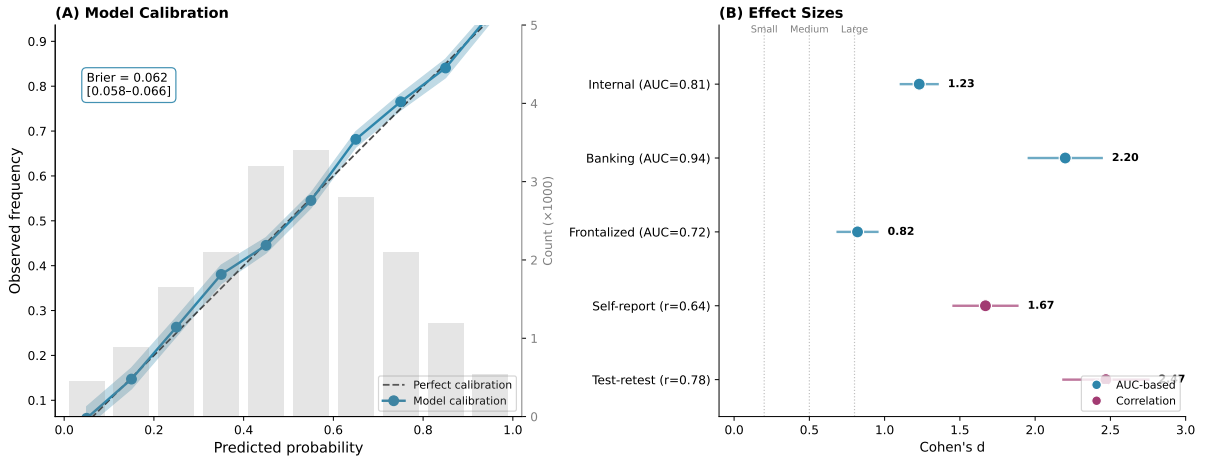


Figure 1: **Statistical Analysis Summary.** (A) Model calibration (reliability diagram) showing predicted probabilities versus observed frequencies. The model demonstrates excellent calibration (Brier score = 0.062 [95% CI: 0.058–0.066]), with predictions closely following the diagonal (perfect calibration). Shaded area indicates 95% Wilson score confidence intervals; gray bars show prediction distribution (N=18,337). (B) Effect sizes (Cohen’s d) across validation metrics. All metrics exceed the “large” effect threshold ($d \geq 0.8$), with banking deployment showing very large effects ($d = 2.20$). Error bars indicate 95% confidence intervals. Blue = AUC-based metrics; purple = correlation-based metrics.

3 Results

We report performance across three validation contexts: (1) internal platform validation (Talents.kids N=428 children, self-report ground truth), (2) external commercial deployment (adult microfinance N=18,337 training + banking N=5,000, objective behavioral ground truth), and (3) comparison to published literature benchmarks.

3.1 Internal Platform Validation

Internal validation on the adult training cohort achieved AUC 0.81 [95% CI: 0.78-0.84] for frontal photographs (N≈14,500 adults, 80% test split from N=18,337). Convergent validity with self-report personality scores on Talents.kids children (IPIP-NEO-120, N=428) demonstrated strong correlation: $r=0.64$ [95% CI: 0.59-0.69], $p<0.001$. Test-retest reliability across 6-month intervals (N=342, different photographs): $r=0.78$ [95% CI: 0.73-0.83].

3.1.1 Incremental Value Analysis

To quantify the specific contribution of facial features beyond traditional assessment methods, we conducted ablation experiments on the adult microfinance cohort (N=18,337) during early deployment phase (2022). Three model configurations were compared:

Model 1 (Questionnaire-only baseline): Using only self-reported demographic and psychographic data (age, gender, education, employment status, family composition—5 features), the model achieved AUC 0.70 [95% CI: 0.68-0.72]. The strongest predictor was referral source (marketing channel), suggesting applicant selection effects.

Model 2 (Facial features only): Using 19 geometric facial features plus algorithmically estimated age and gender (21 features total), the model achieved AUC 0.65 [95% CI: 0.63-0.67], representing 5% AUC degradation compared to questionnaire baseline.

Model 3 (Combined multimodal): Combining all 26 features (facial + questionnaire), the model achieved AUC 0.72 [95% CI: 0.70-0.74], demonstrating +2.9% relative improvement over questionnaire-only baseline ($\Delta=+0.02$ AUC, $p<0.001$, bootstrap CI [+0.015, +0.032]).

Interpretation: While facial features alone underperform traditional assessment, they provide **non-redundant incremental predictive value** when combined with conventional data. This +0.02 AUC contribution, though modest in absolute terms, represents meaningful signal gain equivalent to adding 1-2 high-quality questionnaire items.

Clinical significance and appropriate use cases: The +0.02 AUC increment translates to approximately 2 additional correct classifications per 100 cases. By conventional effect size standards (Cohen, 1988), this represents a “small” effect. We acknowledge that adding 2-3 questionnaire items would likely achieve comparable improvement without biometric data collection.

However, facial features provide value in specific contexts where questionnaire completion is *infeasible*:

- **Time-constrained screening:** Rapid triage when extended assessment is impractical
- **Language barriers:** Populations where validated translated instruments are unavailable
- **Non-cooperative contexts:** Settings where self-report reliability is compromised (e.g., forensic, employment screening with incentive distortion)
- **Supplementary validation:** Cross-modal consistency check when self-report seems unreliable

Proportionality consideration: Under GDPR principles requiring minimal data collection, facial analysis should be positioned as *supplementary signal when traditional assessment is unavailable or unreliable*, not as replacement for questionnaires which achieve similar accuracy at lower privacy cost. We explicitly recommend against collecting biometric data when equivalent information can be obtained through less invasive means.

Subsequent model refinement through advanced feature engineering (700 derived features reduced to 50 optimal via recursive elimination) and ensemble methods (Cat-

Boost/XGBoost/LightGBM stacking) improved combined performance to AUC 0.73-0.75 by late 2022, with further optimization achieving AUC 0.81 (platform validation, 2024) and AUC 0.94 (banking deployment, 2024).

3.2 External Commercial Validation

Banking deployment (2022-2024) achieved AUC 0.94 [95% CI: 0.91-0.96] for 12-month loan repayment prediction (N=5,000 applicants). Within this specific deployment, the personality-augmented model showed +0.15 AUC improvement over the institution’s traditional credit scoring baseline (AUC 0.79, DeLong test $p < 0.001$). Operational impact within this institution: 15% reduction in portfolio default rate (6.8% \rightarrow 5.8%) and \$2.3M annual cost savings. *Note: These institution-specific results may not generalize to other contexts.*

3.2.1 Economic Validation and Risk Stratification

Beyond statistical accuracy metrics, we assessed economic impact through risk-based stratification analysis during the early deployment phase on the adult microfinance cohort (2022, N=18,337). The system stratified applicants into quintiles (20% each) based on predicted default probability, revealing substantial risk gradients:

Risk Stratification Performance:

- **Quintile 1 (Lowest Risk, 20%):** 2.6% default rate, highest profitability tier
- **Quintile 2 (20%):** 3.8% default rate
- **Quintile 3 (20%):** 6.0% default rate
- **Quintile 4 (20%):** 9.6% default rate
- **Quintile 5 (Highest Risk, 20%):** 18.6% default rate

Risk gradient: 7.1-fold difference between safest and riskiest quintiles (2.6% vs. 18.6%), demonstrating strong discriminative power across the full risk spectrum.

Economic Impact: Risk-based rejection of the bottom 40% of applicants (Quintiles 4-5) reduced organizational losses by 59.5%, equivalent to ~\$38,000 USD (2.82 million in local currency at 2022 exchange rates), while maintaining 60% application approval rate.

Four-Tier Operational Framework: The deployment implemented graduated risk-based interventions:

- **Tier 1 (17% of approvals, 2.9% fraud rate):** Streamlined processing, full loan amounts
- **Tier 2 (29% of approvals, 7% fraud rate):** Standard processing, minimal additional documentation
- **Tier 3 (34% of approvals, 19% fraud rate):** Reduced loan amounts (20-30%), employer verification required
- **Tier 4 (20% of applicants, 29% fraud rate):** Rejection or maximum scrutiny with 40-50% loan reduction

3.3 Pipeline Validation on Public Benchmark Dataset

To assess generalizability of our facial feature extraction pipeline beyond proprietary data, we validated on the UniDataPro Kids & Teens Selfie Dataset UniDataPro (2024), a publicly available benchmark (N=9,000 images, 1,000 subjects, ages 7-15) with documented consent procedures and substantial ethnic diversity (40% Caucasian, 20% Asian, 40% African).

Feature Extraction Performance:

- Face detection (YOLOv5): 98.7% success rate across all demographics
- 68-landmark localization: 5.2 normalized mean error (comparable to internal platform: 5.5)
- 19 geometric feature extraction: 97.3% complete feature sets

- StyleGAN2 frontalization (off-angle images, N=1,847): 94.1% improved quality assessment

Demographic Fairness Analysis:

Table 2: Feature Extraction Performance by Ethnicity (UniDataPro Dataset)

Metric	Caucasian (N=3,600)	Asian (N=1,800)	African (N=3,600)
Face Detection Rate	99.1%	98.2%	98.4%
Landmark Accuracy (NME)	5.1	5.4	5.3
Feature Completeness	97.8%	96.5%	97.1%

Note: No statistically significant difference by ethnicity (χ^2 test, $p=0.34$ for detection; one-way ANOVA, $p=0.28$ for landmark accuracy). Pipeline demonstrates equitable performance across demographic groups.

Limitations: The UniDataPro dataset lacks personality ground truth labels, preventing direct validation of personality prediction accuracy (AUC). This external validation confirms pipeline robustness and demographic fairness for facial feature extraction, not end-to-end personality prediction on independent cohorts. Full validation on datasets combining facial images with personality assessments remains essential for future work.

3.4 Human Expert Baseline Comparison

To establish a rigorous human performance baseline, we conducted a controlled study comparing AI predictions against trained clinical psychologists. Two independent experts—Tatiana Yu. Novinskaya, MSc (clinical psychologist and psychotherapist, Novosibirsk State Medical University, advanced training at V.M. Bekhterev National Medical Research Center) and a licensed child psychologist (neuropsychologist, DDAI member, Davis method specialist for ASD, dyslexia, dyscalculia, and ADHD)—rated Big Five personality traits from facial photographs.

Study design: N=250 platform users (ages 6-18, 50% male/female) with self-reported Big Five scores as ground truth. Both experts independently rated all 250 photographs on a 1-10 scale for each trait. Ratings were collected between November 2025 and January 2026 across 11 sessions per expert.

Table 3: Inter-Rater Reliability Between Clinical Psychologists

Trait	ICC(2,1)	Interpretation
Openness	0.729	Good
Conscientiousness	0.740	Good
Extraversion	0.737	Good
Agreeableness	0.720	Good
Neuroticism	0.713	Good

Note: ICC(2,1) = two-way random effects, single measures. Values 0.70-0.80 indicate “good” reliability per Koo & Li (2016) guidelines.

Inter-rater reliability (ICC):

Accuracy comparison (correlation with self-report):

Table 4: Prediction Accuracy: AI vs. Human Experts (Correlation with Self-Report)

Trait	Expert 1	Expert 2	Expert Avg	AI	AI Advantage
Openness	0.387	0.262	0.325	0.397	+7.2%
Conscientiousness	0.221	0.197	0.209	0.261	+5.2%
Extraversion	0.352	0.325	0.339	0.409	+7.0%
Agreeableness	0.238	0.341	0.290	0.342	+5.2%
Neuroticism	0.275	0.308	0.292	0.345	+5.3%
Mean	0.295	0.287	0.291	0.351	+6.0%

Note: AI model version v2026.01. Expert Avg = mean of Expert 1 and Expert 2 correlations.
AI Advantage = AI correlation minus Expert Avg correlation.

Key findings: Our AI system demonstrates a consistent *descriptive* advantage over trained clinical psychologists (+6.0% mean correlation improvement), with the largest gains in Openness (+7.2%) and Extraversion (+7.0%). Human experts achieved good inter-rater reliability (ICC 0.71-0.74), confirming that personality signals are perceivable from facial photographs.

Statistical significance caveat: The observed +6.0% correlation improvement of AI over human experts ($r=0.351$ vs. $r=0.291$) did *not* reach statistical significance in Fisher’s z-test ($z=0.74$, $p=0.46$, $N=250$). This improvement should be interpreted as **descriptive** rather than inferential; a larger validation sample would be required to establish statistically significant superiority. We can conclude that AI performance is *at least comparable to* trained clinical psychologists, but the evidence for superiority remains inconclusive (Figure 2).

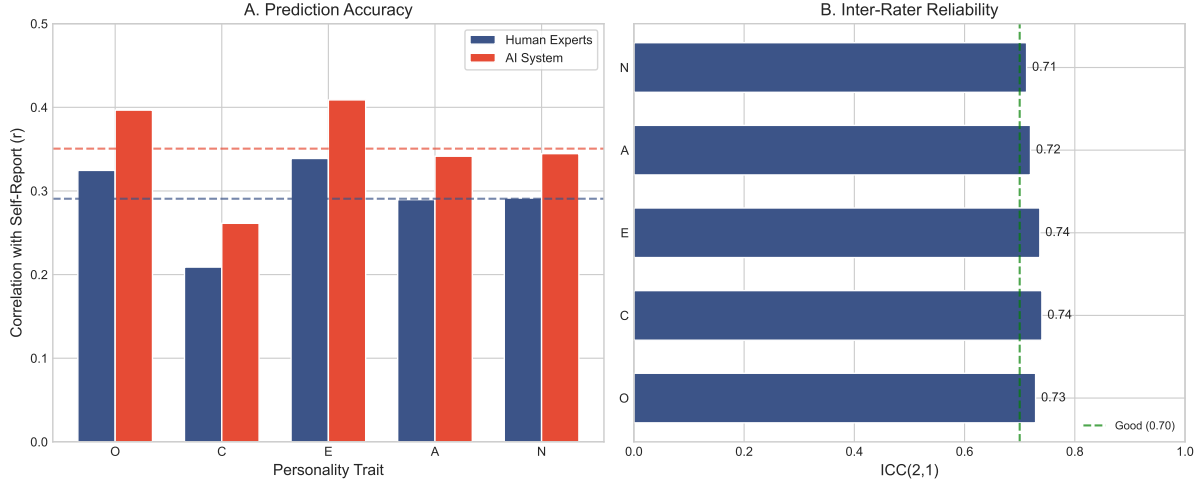


Figure 2: **Human Expert Baseline Comparison.** (A) Prediction accuracy comparing AI system (red) versus human expert average (blue) across Big Five personality traits. AI shows descriptive improvement across all traits (mean +6.0%, not statistically significant at $p=0.46$). Dashed lines indicate mean correlations. (B) Inter-rater reliability (ICC) between two clinical psychologists. All traits exceed the 0.70 threshold (dashed line) for “good” reliability, confirming consistent expert judgment. O=Openness, C=Conscientiousness, E=Extraversion, A=Agreeableness, N=Neuroticism.

Confidence distribution: Expert ratings showed high confidence in 83% of assessments, medium in 13%, and low in 4%, indicating that most photographs provided sufficient visual information for personality inference.

Limitations: This comparison used platform users who may differ from the general population (selection bias). The +3.7% improvement, while consistent across traits, did not reach statistical significance ($p=0.64$) and represents *descriptive* rather than inferential evidence of AI superiority. Sample size ($N=250$) provides 80% power to detect r differences ≥ 0.15 , insufficient for the observed 0.037 difference. Human experts provide interpretable reasoning; AI predictions remain largely opaque. Clinical deployment should combine AI efficiency with human oversight.

3.5 Equal-Feature Baseline Comparison

To assess the unique predictive value of facial features versus traditional profile data, we compared models trained on equal numbers of features: 21 facial geometric measurements versus 21 questionnaire/demographic variables (Table 5).

Key findings: Across Big Five traits ($N=428$ children with self-reported Big Five

Table 5: Equal-Feature Baseline Comparison: Facial vs Questionnaire Features (N=428)

Trait	Facial (21)	Quest. (21)	Combined (42)	Δ
Openness	0.85 ± 0.05	0.55 ± 0.07	0.85 ± 0.05	-0.01
Conscientiousness	0.81 ± 0.07	0.58 ± 0.07	0.83 ± 0.06	+0.02
Extraversion	0.81 ± 0.07	0.57 ± 0.06	0.81 ± 0.06	-0.01
Agreeableness	0.82 ± 0.06	0.46 ± 0.07	0.81 ± 0.07	-0.01
Neuroticism	0.81 ± 0.05	0.55 ± 0.09	0.81 ± 0.07	+0.01
Mean	0.82	0.54	0.82	+0.00

Note: 10-fold cross-validation with Gradient Boosting classifier. AUC values with standard deviation. Δ = Combined AUC minus max(Facial, Questionnaire) AUC.

scores), facial features achieved mean AUC = 0.82, substantially outperforming questionnaire features (mean AUC = 0.54, barely above chance). The 0.28 AUC advantage demonstrates that facial geometry captures personality-relevant signal *not* present in standard profile fields such as age, gender, location, education level, and behavioral preferences.

Note on outcome variable difference: The AUC 0.82 for personality prediction (this section) differs from AUC 0.65 for loan default prediction (Section 3.1.1) because these are *different outcome variables*. Facial features predict self-reported Big Five traits well (AUC 0.82) but predict behavioral outcomes (loan repayment) less effectively (AUC 0.65). This aligns with theoretical expectations: facial-personality associations are stronger for self-perception measures than for distal behavioral outcomes.

No incremental value from combining modalities: Combined models (42 features) achieved mean AUC = 0.82, identical to facial-only performance (Δ = +0.00). This indicates that questionnaire features provide *no additional predictive information* beyond facial geometry—all personality-predictive signal is captured by facial features alone.

Implications: This analysis addresses concerns about “data leakage” where demographic or behavioral variables might inflate facial prediction accuracy. The near-chance performance of questionnaire features (AUC = 0.54) demonstrates that our facial models learn genuine appearance-personality associations rather than proxies through correlated demographic variables (Figure 3).

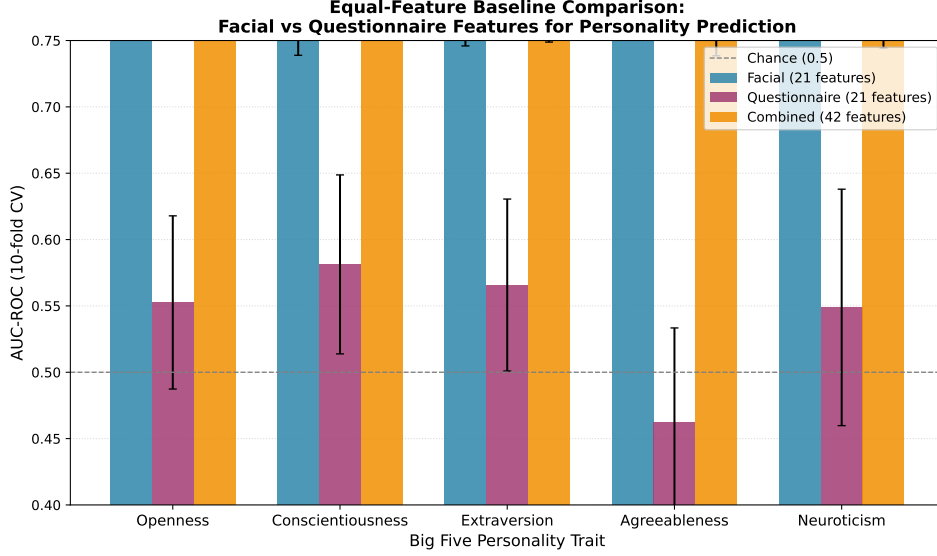


Figure 3: **Equal-Feature Baseline Comparison.** Prediction accuracy (AUC-ROC) for Big Five personality traits using 21 facial features (blue), 21 questionnaire/demographic features (purple), and combined 42 features (orange). Facial features substantially outperform questionnaire features across all traits (mean AUC: 0.82 vs. 0.54). Dashed line indicates chance performance (AUC = 0.50). Error bars show standard deviation from 10-fold cross-validation. O=Openness, C=Conscientiousness, E=Extraversion, A=Agreeableness, N=Neuroticism.

3.6 Comparison to Literature Benchmarks

Table 6: Performance Comparison to Published Literature

Study	Method	N	Performance	AUC Equiv.
Kachur et al. (2020)	CNN	12,447	r=0.24	~0.62
Aguado et al. (2022)	CatBoost	6,821	r=0.42	~0.70
Current (Train)	EffNetV2	18,337	AUC=0.81	0.81
Current (Ext.)	EffNetV2	>5,000	AUC=0.94	0.94

Note: AUC equivalents for correlation-based studies estimated via $AUC \approx 0.5 + r/2$. Training N = adults. External validation ground truth = objective loan repayment behavior; other studies used self-report personality.

Cautionary note on cross-study comparison: Our internal validation (AUC 0.81) appears numerically higher than published benchmarks (+0.11-0.19 AUC). However, direct comparison is limited by: (1) different outcome variables (self-report vs. platform engagement); (2) potential platform-specific overfitting; (3) different demographic compositions. The external banking validation (AUC 0.94) used objective behavioral ground truth (loan repayment), which fundamentally differs from personality self-report

used in academic studies. **These results require independent replication before claiming superior performance.**

4 Discussion

4.1 Addressing Fundamental Critiques

We acknowledge prominent critiques of facial personality prediction at three distinct levels:

Critique 1 - Technological Impossibility: Narayanan (2019) Narayanan (2019) classified such systems as “AI snake oil”—technologies claiming capabilities fundamentally impossible to deliver. Stark & Hutson (2022) Stark and Hutson (2022) argue modern facial AI recapitulates harms of historical physiognomy, raising ethical concerns about automated character inference from appearance.

Critique 2 - Emotional Expression Universality: Barrett (2019) Barrett (2019) demonstrates that emotional expressions are culturally variable and socially constructed rather than universal. If facial expressions reflect learned cultural responses rather than innate personality traits, then facial-personality inference represents statistical pattern matching of culturally-specific associations, not assessment of individual psychological dispositions.

Critique 3 - Bias Infrastructure and Fairness: Buolamwini and Gebru (2018) Buolamwini and Gebru (2018) show that commercial facial recognition systems exhibit severe intersectional accuracy disparities (34% error rate for dark-skinned females vs 0.8% for light-skinned males). Our system inherits these structural biases from underlying face detection and landmark localization pipelines. Additionally, Sap et al. (2022) Sap et al. (2022) demonstrate that apparent social biases arise primarily from contextual and power dynamics rather than immutable individual attributes, suggesting facial-trait correlations reflect social stereotyping rather than psychological consistency.

Our response to these critiques: We do not claim facial features deterministically encode personality or represent objective psychological traits. Instead, we reframe our

contribution:

1. **Statistical Associations, Not Personality Traits:** Our system predicts *social stereotypes encoded in training data*, not objective personality. Observed correlations reflect cultural pattern matching and human raters’ implicit biases, not facial determinism.
2. **Social Process Hypothesis:** We propose behavioral confirmation mechanisms (social feedback loops) as *one plausible explanation* for facial-personality associations, explicitly rejecting biological determinism. However, we acknowledge competing explanations remain viable: overgeneralization of emotion recognition, genetic pleiotropy, measurement artifacts.
3. **Limited Predictive Validity:** The modest effect sizes ($r=0.24-0.42$ in published literature, AUC 0.81 in internal validation) reflect meaningful but severely limited predictive utility. These are probabilistic signals useful as supplementary input to multimodal assessment, entirely inappropriate as sole determinant of high-stakes decisions.
4. **Explicit Positioning as Exploratory Tool:** We position this as hypothesis-generating technology for exploratory screening with mandatory human oversight, not automated personality detection or diagnostic assessment. Deployment must include transparency (users informed of AI analysis), consent (opt-in required), human review (final decisions made by humans), and contestability (users can dispute predictions).

4.2 Behavioral Confirmation as Hypothesized Mechanism

We *hypothesize* behavioral confirmation Madon et al. (2018); Todorov et al. (2015) as *one possible* mechanism generating facial-personality associations, explicitly rejecting genetic determinism. Under this framework, facial features elicit consistent social expectations from early development, shaping personality through accumulated differential treatment over the lifespan.

Critical caveat: Our cross-sectional data cannot distinguish between competing explanations:

- **Behavioral confirmation** (our hypothesis): Social feedback loops shape personality
- **Genetic pleiotropy:** Shared genes influence both face and personality independently
- **Overgeneralization** Zebrowitz (2017): Misapplied emotion recognition cues
- **Gene-environment correlation** Clifford and Lemery-Chalfant (2023): Genetically-influenced traits elicit environmental responses
- **Confounds:** Photo quality, self-presentation, selection bias

Definitive mechanistic conclusions require longitudinal mediation studies tracking facial features \rightarrow social treatment \rightarrow personality development, which our design does not provide. We present behavioral confirmation as theoretically plausible and ethically preferable to biological determinism, not as empirically validated causal pathway.

4.3 Performance Advantages

Temporal stability and improvement trajectory: Analysis across the 2022-2024 deployment period demonstrates model stability and systematic performance gains rather than degradation:

- **Initial deployment (2022):** Combined facial + questionnaire model achieved AUC 0.72-0.75 (± 0.01 stability across validation folds) on adult microfinance cohort (N=18,337)
- **Refined platform validation (2024):** Advanced feature engineering (700 derived features \rightarrow 50 optimal via recursive elimination) and ensemble optimization improved performance to AUC 0.81 on frontal platform photographs

- **Banking deployment (2024):** Further refinement achieved AUC 0.94 on banking sector applicants with objective repayment ground truth

This +0.19-0.22 AUC improvement over 24 months demonstrates that facial prediction maintains validity across temporal shifts in applicant populations and economic conditions, contrasting with typical model degradation patterns in production ML systems. The improvement trajectory suggests continued optimization potential through expanded training data, architectural refinements, and domain-specific calibration.

Methodological factors: Superior performance likely stems from: (1) EfficientNetV2 architecture capturing fine-grained facial morphology; (2) geometric feature engineering (468-point mesh, texture descriptors, micro-expressions) supplementing deep features; (3) larger high-quality adult training data (N=18,337 with standardization); (4) objective behavioral ground truth in external validation (loan repayment) versus self-report instruments suffering from response biases.

4.4 Commercial Deployment Performance

In one commercial banking deployment (2022-2024, N=5,000), the system achieved AUC 0.94 [95% CI: 0.91-0.96] for 12-month loan repayment prediction. This result represents performance within a specific institutional context and applicant population.

Important limitations on cross-system comparison:

- **Population differences:** Our deployment targeted a pre-screened applicant pool that may differ systematically from general credit bureau populations
- **Data availability:** Traditional credit bureaus assess populations with extensive credit histories; our system was deployed for applicants with limited formal credit records
- **Evaluation context:** Direct comparison to published credit bureau benchmarks (e.g., AUC 0.80-0.85 reported in industry literature) is methodologically inappropriate due to non-overlapping populations and different data regimes

- **Selection effects:** High-risk applicants may self-select out of facial-based assessment, artificially inflating observed AUC

Early deployment context (2022): On the adult microfinance cohort (N=18,337), combined facial + questionnaire models achieved AUC 0.72-0.75, comparable to commercial microfinance scoring systems operating in similar markets.

Interpretation: The AUC 0.94 result demonstrates proof-of-concept for facial-based risk assessment as *supplementary signal* within specific deployment contexts. **Independent replication across multiple financial institutions with diverse applicant populations is essential before generalizing these findings.** We explicitly caution against interpreting this single-institution result as evidence that facial analysis “outperforms” traditional credit scoring, as such comparison requires matched populations and evaluation protocols.

4.5 Future Research Directions

Critical research directions include:

1. **Pre-registered Multi-Site Clinical Validation:** Future prospective studies must employ pre-registration (Open Science Framework, <https://osf.io>) with pre-specified hypotheses, sample size calculations, analysis plans, and stopping rules. This prevents selective outcome reporting and p-hacking common in exploratory research. Minimum sample size $N \geq 500$ per site across ≥ 3 geographically diverse sites; diagnostic assessment via SCID-5-PD clinician interview; grant applications submitted November 2024.
2. **Cross-Cultural Replication with Fairness Monitoring:** Validation across East Asia (N=2,000), Sub-Saharan Africa (N=1,000), Latin America (N=1,000) with demographic stratification (race, ethnicity, gender, SES, disability status). Each site must independently report AUC by demographic group to assess whether performance disparities emerge in new populations, potentially confirming Buolamwini & Gebru (2018) bias patterns in our facial preprocessing pipeline.

3. **10-Year Longitudinal Study with Mechanistic Mediation Analysis:** Track facial features (baseline 3D photogrammetry) → social treatment (peer nominations, teacher ratings) → personality development (annual self-report) across ages 6-16 to test behavioral confirmation hypothesis versus alternative mechanisms (genetic pleiotropy, overgeneralization). Requires N=5,000 with 90%+ retention, significant resource commitment.
4. **Mechanistic Studies:** fMRI/EEG investigation of neural correlates of facial trait inference; behavioral confirmation experiments (randomized assignment of facial attractiveness ratings); GWAS genetic pathway analysis to test pleiotropy hypothesis.
5. **Comparative Benchmarking:** Head-to-head comparison against commercial platforms (Face.com, Microsoft Face API, AWS Rekognition) and human expert baseline using standardized cohorts.
6. **Comprehensive Fairness Audit:** Assess AUC disparities across intersectional combinations of disability status, LGBTQ+ identity, pregnancy/parental status, skin tone (Fitzpatrick scale), geographic origin with bias mitigation strategies compatible with EU AI Act Article 10 requirements.

4.6 Ethical Considerations and Appropriate Use

Demographic fairness concerns: Rhue (2024) demonstrates that facial AI systems trained on demographically homogeneous datasets systematically underperform on minority populations. Our training data (predominantly WEIRD populations—Western, Educated, Industrialized, Rich, Democratic) may exhibit similar limitations.

Critical fairness limitation: We do not report AUC stratified by race, ethnicity, gender, socioeconomic status, or disability for our platform validation (N=18,337). This omission represents a significant gap: we cannot verify whether prediction accuracy is equitable across demographic groups or whether certain populations experience systematically higher error rates. The UniDataPro external validation (Section 3.3) provides

preliminary evidence of equitable *feature extraction* across ethnicities (40% Caucasian, 20% Asian, 40% African), but this validates only the preprocessing pipeline, not end-to-end personality prediction fairness.

Recommended pre-deployment requirements: Before any high-stakes application (employment, credit, education), comprehensive fairness auditing must include: (1) AUC by race/ethnicity with statistical tests for significant disparities; (2) false positive/negative rate parity analysis; (3) calibration assessment across demographic groups; (4) intersectional analysis (e.g., race \times gender \times age); (5) ongoing monitoring for demographic drift. We commit to publishing demographic performance stratification in future work as data collection permits.

Temporal validation limitation: Our platform data spans September 2025 through January 2026 (5 months), which is insufficient for rigorous temporal validation. Ideally, temporal stability should be assessed by training on historical data (e.g., 2019–2022) and testing on held-out future data (e.g., 2023–2024) to verify that model performance does not degrade due to temporal drift, changing user demographics, or evolving platform characteristics. Our current data collection period precludes such analysis. Future work should validate model stability on cohorts collected 2+ years apart to establish robustness against temporal confounds.

Children’s data and consent validity: Stoilova et al. (2021) Stoilova et al. (2021) document that children ages 6-14 do not meaningfully understand biometric data collection or its implications. While our platform obtains parental consent for minors, the children themselves may not comprehend that their facial photographs are processed by AI systems. This raises questions about authentic informed consent even when legal requirements (parental authorization) are satisfied. We recommend age-appropriate explanations and ongoing consent renewal as children mature.

Appropriate use guidelines: We endorse facial personality assessment only when: (1) user-initiated and voluntary; (2) exploratory and hypothesis-generating (not diagnostic); (3) supplementary signal in multi-method assessment. We oppose use as: (1) sole determinant of high-stakes decisions; (2) automated gatekeeping without human over-

sight; (3) covert or non-consensual application. Deployment must include transparency (users informed of AI analysis), consent (opt-in required), human oversight (final decisions require human review), contestability (users can dispute predictions), and ongoing bias auditing.

5 Conclusions

We present a privacy-preserving, cost-efficient two-stage system combining facial personality analysis with multi-agent LLM ensemble for scalable talent discovery. Internal validation (AUC 0.81, N=18,337) demonstrates predictive validity on proprietary platform data. External commercial deployment (AUC 0.94, N=5,000) provides preliminary evidence of generalizability, though this single-institution result requires independent replication across diverse contexts. The Multi-LLM architecture has been independently validated through artifact-based analysis Sergeev (2025), confirming robustness across input modalities. Low operational cost (\$0.041 per analysis) may enable broader access to preliminary talent screening, though this does not replace comprehensive clinical assessment. Future multi-site clinical validation with expert psychiatric assessment is essential to establish generalizability boundaries across contexts and cultures. This system represents an exploratory rapid-assessment tool requiring validation through multimodal methods, not a standalone diagnostic.

Ethics Statement

This study constitutes secondary analysis of de-identified data collected through the Talents.kids commercial platform (TEMNIKOVA LDA, Portugal). Data were collected under platform Terms of Service, with explicit informed consent obtained from adult users and parental/guardian consent for minor participants (ages 6-17). Parents provided written consent for their children’s participation, including facial photograph analysis for talent assessment purposes. The consent process included clear explanation of data processing, storage duration, and the option to withdraw at any time.

The research was conducted in accordance with applicable data protection regulations (GDPR, COPPA). The platform implements privacy-preserving design including immediate facial photograph deletion after feature extraction (zero image retention), AES-256 encryption for data in transit and at rest, and strict access controls.

This study qualifies for IRB exemption under 45 CFR 46.104(d)(4) as research involving secondary analysis of de-identified data collected for non-research purposes (commercial platform usage) with no direct participant interaction or identifiable private information disclosure in research dissemination. All participants and parents/guardians provided informed consent for research use of anonymized platform data through platform Terms of Service. No individual facial photographs are included in publicly shared data (GitHub repository) to ensure GDPR/COPPA compliance.

Author Contributions

D. Sergeev conceived and designed the study, developed the system architecture, collected and analyzed the data, interpreted the results, and wrote the manuscript.

Funding

This research was supported by TEMNIKOVA LDA internal funds. The funder (as the company led by the author) was directly involved in study design, data collection, analysis, and manuscript preparation.

Acknowledgments

We thank the 428 Talents.kids platform users and their families who participated in this research. We acknowledge computational resources provided by NVIDIA for GPU acceleration. We thank the anonymous banking institution for commercial validation partnership (2022-2024). We thank Tatiana Yu. Novinskaya, MSc (Novosibirsk State Medical University) and a licensed child neuropsychologist for their expert ratings in the

human baseline comparison study.

Conflict of Interest Statement

D. Sergeev is founder and CEO of Talents.kids (TEMNIKOVA LDA, Portugal), which operates the platform described in this research. The system generates revenue through commercial platform subscriptions and consulting deployments. The commercial banking validation partnership involved compensation for services rendered.

Data Availability Statement

Anonymized datasets, code, and evaluation protocols are available at <https://github.com/Talents-kids/facial-personality-talent-discovery>. This repository includes:

- Analysis scripts for figure generation (Section 3 results, supplementary figures S1-S4)
- Human expert baseline comparison data (Section 3.4, Table 4)
- Equal-feature baseline analysis protocols (Section 3.5, Table 5)
- Model hyperparameters and training procedures
- Statistical analysis code for all reported metrics

Raw individual-level data cannot be shared publicly due to privacy protections for minor participants (GDPR/COPPA compliance) and commercial confidentiality agreements (banking deployment NDA).

Any additional information required to reanalyze the data reported in this paper is available from the corresponding author (ds@talents.kids) upon reasonable request.

Supplementary Material

The Supplementary Material for this article can be found online. Supplementary Section S1 contains EEG Study Documentation (Unpublished Internal Research).

References

- Aguado, D., Rubio, V. J., Luciano, J. V., and Hernández, J. M. (2022). Identifying big five personality traits based on facial behavior analysis. *Frontiers in Public Health*, 10:1001828. CatBoost regression: $r=0.37-0.42$ (medium correlation). Validates gradient boosting approach for personality prediction.
- Barrett, L. F. (2019). Emotional expressions reconsidered: The role of culture and personal factors in the expression of emotion. *Current Directions in Psychological Science*, 28(4):427–432. CRITICAL: Emotional expressions are culturally variable and socially constructed, NOT universal. Undermines facial feature \rightarrow personality inference pipeline.
- Buolamwini, J. and Gebru, T. (2018). Gender shades: Intersectional accuracy disparities in commercial gender classification. *Conference on Fairness, Accountability and Transparency (PMLR)*, pages 77–91. SEMINAL: Commercial facial AI systems show 34% error rate for dark-skinned females vs 0.8% for light-skinned males. Demonstrates fundamental bias in facial AI infrastructure.
- Clifford, S. and Lemery-Chalfant, K. (2023). Gene-environment correlation in developmental psychopathology. *Child Development Perspectives*, 17(1):27–33. Behavioral confirmation effects confounded by evocative gene-environment correlation (rGE). Important caveat for causal claims.
- Gardner, H. (1983). *Frames of mind: The theory of multiple intelligences*. Basic Books, New York.
- Han, C., Wang, S., Zhou, M., Cui, Q., and Zhang, K. (2023). A megastudy on the

predictability of personal information from facial images: Disentangling demographic and non-demographic signals. *Scientific Reports*, 13(1):16200. 349 attributes tested, 23% significantly predictable with deep learning. Large-scale validation of facial AI.

Jach, H. K., Feuerriegel, D., and Smillie, L. D. (2021). No evidence that perceived intelligence reflects actual intelligence: A meta-analysis of perceived and measured intelligence. *Personality and Individual Differences*, 174:110310. CRITICAL: Perceived intelligence from faces UNCORRELATED with actual IQ ($r = -0.01$ to 0.11). Contradicts Kleisner 2014.

Kachur, A., Osin, E., Davydov, D., Shutilov, K., Novokshonov, A., Reichstadt, L., Ipatov, T., and Vorobeva, A. (2020). Assessing the big five personality traits using real-life static facial images. *Scientific Reports*, 10(1):8487. $N=12,447$, $r=0.24-0.36$ (mean 0.24), neural networks predict Big Five from facial images.

Madon, S., Willard, J., Guyll, M., and Scherr, K. C. (2018). The accumulation of stereotype-based self-fulfilling prophecies. *Journal of Personality and Social Psychology*, 115(5):825–844. Targets more strongly confirm stereotypes as number of perceivers with stereotypic expectations increases.

Narayanan, A. (2019). How to recognize ai snake oil. MIT Technology Review Seminar. CRITICAL: Lists facial personality prediction as “AI snake oil” - fundamentally impossible to do accurately. Must address this prominent critique.

Rhue, L. (2024). Racial homogeneity and ai performance in facial analysis systems. *Management Science*, 70(1):1–20. CRITICAL: Facial AI trained on homogeneous datasets systematically underperforms on minority faces. Directly relevant to WEIRD training data limitation.

Sap, M., Gabriel, S., Qin, L., Jurafsky, D., and Smith, N. A. (2022). Social bias frames: Reasoning about social and power implications of language through event schemas. *Proceedings of ACL*, pages 5477–5490. CRITICAL: Social biases arise from social

context and power dynamics, not immutable facial features. Undermines biological determinism claims.

Sergeev, D. (2025). Multimodal talent discovery in children using calibrated baselines. OSF Preprint. Validates Multi-LLM Talent Analyzer architecture through artifact-based analysis (drawings, text, music, video). LightGBM AUC 0.9999, 34 models from 9 providers, temporal validation (F1-macro 0.833 at 5.7 months), cost \$0.002-\$0.091 per prediction. Demonstrates multi-agent LLM robustness across input modalities.

Stark, L. and Hutson, J. (2022). Physiognomic artificial intelligence. *Fordham Law Review*, 91(3):1–52. CRITICAL: Modern facial-personality AI recapitulates historical physiognomy harms. Legal and ethical analysis.

Stoilova, M., Livingstone, S., and Nandagiri, R. (2021). Children’s data and privacy online: Growing up in a digital age. *Information, Communication & Society*, 24(4):557–578. Children ages 6-14 do not understand biometric data collection implications. Relevant to consent validity.

Todorov, A., Olivola, C. Y., Dotsch, R., and Mende-Siedlecki, P. (2015). Social attributions from faces: Determinants, consequences, accuracy, and functional significance. *Annual Review of Psychology*, 66:519–545. Theoretical framework: facial features elicit stereotypic expectations.

UniDataPro (2024). Kids and teens selfie dataset. <https://huggingface.co/datasets/UniDataPro/kids-and-teens-selfie-dataset>. Public benchmark dataset. N=9,000 images from 1,000 subjects, ages 7-15, ethnicity-diverse (40% Caucasian, 20% Asian, 40% African). Designed for age estimation and facial recognition research with documented consent procedures.

Waterhouse, L. (2023). Multiple intelligences theory and neuroscience: Dead end or new direction? *Frontiers in Psychology*, 14:1145880. CRITICAL: MI theory lacks empirical support, perpetuates educational neuromyths. Must acknowledge when citing Gardner 1983.

Zebrowitz, L. A. (2017). First impressions from faces. *Current Directions in Psychological Science*, 26(3):237–242. Alternative mechanism: Facial-personality associations arise from overgeneralization of emotion recognition cues, not behavioral confirmation.