# EZ-Sort: Efficient Pairwise Comparison via Zero-Shot CLIP-Based Pre-Ordering and Human-in-the-Loop Sorting

Yujin Park
Hanyang University
Seoul, Republic of Korea
yujin1019a@hanyang.ac.kr

Haejun Chung[*]
Hanyang University
Seoul, Republic of Korea
haejun@hanyang.ac.kr

Ikbeom Jang[*]
Hankuk University of Foreign Studies
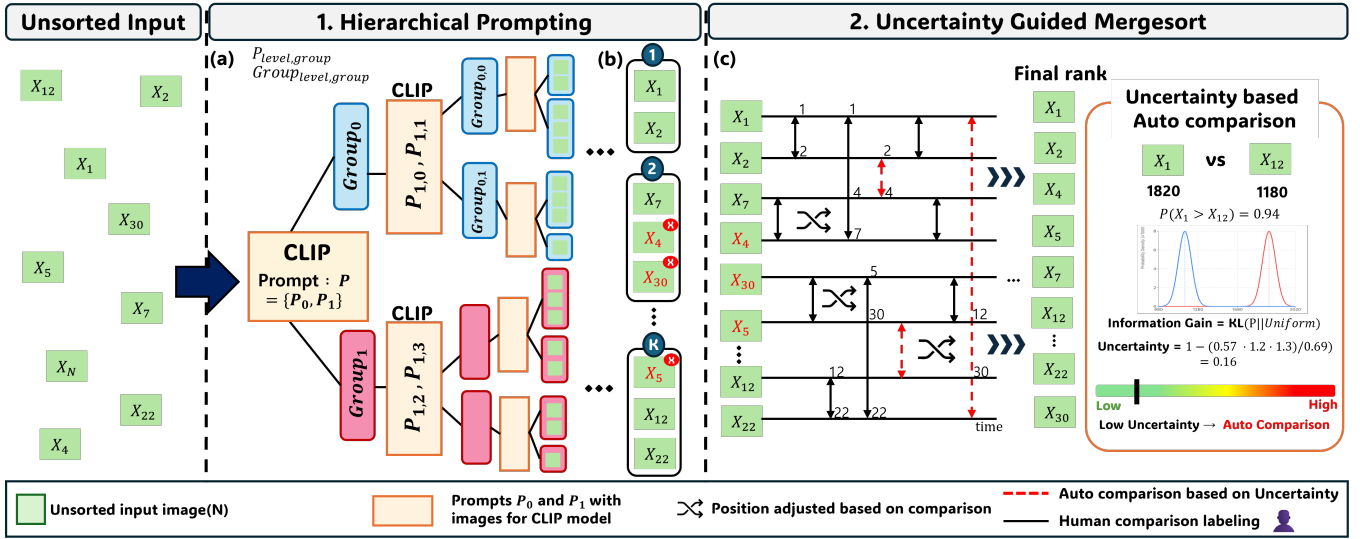Yongin, Republic of Korea
ijang@hufs.ac.kr

Figure 1: Overview of the EZ-Sort framework. The framework operates in three stages: (a) CLIP-based zero-shot *pre-ordering* is performed via hierarchical prompting that recursively groups unsorted images using binary prompts. (b) The resulting fine-grained groups are merged into $k$ coarse buckets, and each image is assigned an Elo score based on its bucket ID and CLIP confidence. (c) An uncertainty-aware MergeSort selectively routes high-uncertainty comparisons to human annotators while automatically resolving confident ones.

## Abstract

Pairwise comparison is often favored over absolute rating or ordinal classification in subjective or difficult annotation tasks due to its improved reliability; however, exhaustive comparisons require a massive number of annotations ($O(n^2)$). Recent work [4] greatly reduced the annotation burden ($O(n \log n)$) by actively sampling pairwise comparisons using a sorting algorithm. We further improve annotation efficiency by 1) roughly pre-ordering items using the Contrastive Language-Image Pre-training (CLIP) model hierarchically without training and 2) replacing easy, obvious human comparisons with automated comparisons. The proposed EZ-Sort first produces a CLIP-based zero-shot *pre-ordering*, then initializes bucket-aware Elo scores, and finally runs an uncertainty-guided

human-in-the-loop MergeSort. Validation was conducted using various datasets: face-age estimation (FGNET) [21], historical image chronology (DHCI) [20], and retinal image quality assessment (EyePACS) [19]. It showed that EZ-Sort reduced human annotation cost by 90.5% compared to exhaustive pairwise comparisons and by 19.8% compared to prior work [4] (when $n = 100$) while improving or maintaining inter-rater reliability. These results demonstrate that combining CLIP-based priors with uncertainty-aware sampling yields an efficient and scalable solution for pairwise ranking. Code available at https://github.com/yujinPark02/EZ-Sort-CIKM2025.

## CCS Concepts

• **Information systems** → *Data labeling*; • **Computing methodologies** → *Ranking*.

## Keywords

Pairwise comparison, Human-in-the-loop sorting, VLM-based pre-ordering, Annotation, Labeling

---

## 1 Introduction

Pairwise comparison is widely preferred in subjective annotation tasks, including perceptual quality assessment, face-age estimation, and medical image triage [37], due to its superior inter-rater reliability compared to absolute or ordinal ratings. However, exhaustive pairwise labeling incurs a quadratic annotation burden $O(n^2)$, quickly becoming infeasible as dataset sizes grow. This annotation bottleneck has tangible implications for large-scale clinical diagnostics, population health studies, and the preservation of historical records, where subject-matter expertise is scarce and annotation budgets are limited. This scalability challenge can be addressed by leveraging computational sorting algorithms, such as MergeSort.

However, traditional methods, such as the Bradley–Terry–Luce model [5] and active learning strategies [24, 31], typically assume uniform priors and overlook opportunities to leverage the existing semantic structure in the data. More recently, sorting-based approaches have incorporated active query selection to reduce the number of comparisons [4].

To further alleviate the annotation burden, we propose incorporating vision-language models (VLM), such as CLIP [13], to provide a strong starting point for the sorting process. Because such models are pre-trained on hundreds of millions of image-text pairs, proper prompts enable a coarse initial ranking of given items to be labeled. This leads to a significant reduction in the number of comparisons needed for sorting. We iteratively execute this process hierarchically to improve accuracy. We also propose replacing human comparisons with automated comparisons for item pairs with low uncertainty. By combining this prior knowledge with an uncertainty-guided comparison selection strategy, we aim to provide a complementary perspective to existing methods, focusing on the efficient utilization of both model priors and human expertise.

Our key contributions include (1) leveraging pre-trained vision-language priors to reduce the initial annotation search space significantly, (2) applying this VLM-based pre-ordering hierarchically for improved accuracy, and (3) introducing a novel uncertainty-guided sorting strategy to prioritize human annotation resources intelligently.

## 2 Methods

EZ-Sort follows a three-stage pipeline: (i) hierarchical CLIP prompting provides a zero-shot semantic pre-ordering (i.e., ordinal classification) of images, (ii) bucket-aware Elo scores initialize priors, and (iii) Kullback-Leibler(KL) based MergeSort routes only uncertain pairs to annotators. [32, 35].

This pipeline reflects a dual-process structure, with automatic, low-uncertainty comparisons (System 1) and human deliberation for high-uncertainty cases (System 2). The subsequent subsections provide a detailed overview of each stage.
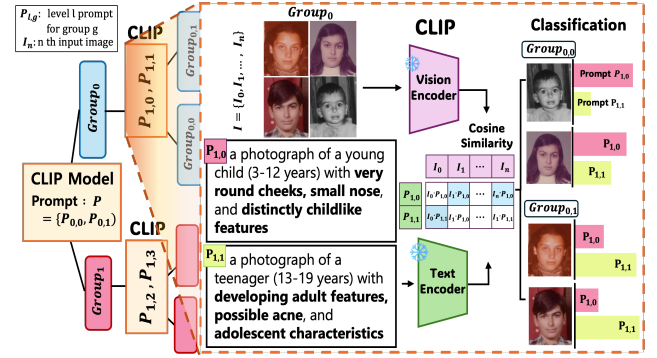


**Figure 2: CLIP-based hierarchical classification at Level 1. Input images are compared against two age-related prompts (child vs. teenager) using CLIP cosine similarity and assigned to subgroups based on the highest score.**

### 2.1 Hierarchical Prompt Design for Zero-Shot Classification

Traditional classification often struggles with ambiguous subjective attributes. To address this, we introduce a hierarchical prompting strategy that exploits CLIP's zero-shot capabilities through multi-level binary decisions. Inspired by how machine learning has optimized sorting algorithms [1, 8, 25], our method enhances classification robustness by decomposing complex decisions into more tractable binary ones.

To evaluate its benefit over single-level prompting, we tested two flat baselines: a minimal template $P_{\text{simple}} = \{$"a photograph of a $c$"$\}$ and an enhanced variant incorporating GPT-4-generated attributes $P_{\text{flat-detailed}} = \{$"a photograph of $c$ with $\mathcal{A}_c$"$\}$. Experiments across three datasets demonstrate that hierarchical prompts outperform both baselines, yielding an improvement of up to 2.0 MAE. This performance gain arises from (1) **decision decomposition**, which replaces $k$-way classification with $\log_2 k$ binary steps where CLIP excels, and (2) **progressive refinement**, which focuses each decision on discriminative features at that granularity. By mimicking human coarse-to-fine reasoning, our hierarchical strategy enhances both interpretability and classification accuracy, particularly in visually ambiguous domains.

*2.1.1 Prompt Generation Methodology.* We replace flat $k$-class classification with an adaptive hierarchical binary structure, dynamically yielding groups until no meaningful visual distinctions remain. At each level $\ell$, we define $P_\ell = \{p_{\ell,0}, p_{\ell,1}, \ldots, p_{\ell,2^\ell-1}\}$, building upon recent advances in prompt learning [2, 18].

> ***Automated Hierarchical Prompt Generation:***
> *"Given the domain **[domain]** with range **[range]**, iteratively divide each group into two visually distinguishable sub-groups using observable anatomical or textural features. Continue dividing until no further meaningful visual distinctions are identifiable. **Avoid** behavioural or contextual clues."*

For face-age estimation (*domain = face, range = 0–60+ years*) this produces prompts such as *"rounded cheeks, large forehead"* (infants)

versus *"defined cheekbones, mature jawline"* (adults); the recursion stops once visual distinctions become ambiguous.

**Adaptive Depth Criterion.** Depth is not fixed; generation halts when GPT-4 signals that further splits are unreliable, resulting in 3 to 5 levels in practice, balancing granularity with annotation load.

**Group Assignment.** Binary outcomes $c_{i,\ell} \in \{0, 1\}$ are combined into a final group index

$$g_i = \sum_{\ell=1}^{d_i} c_{i,\ell} \, 2^{\ell-1}, \qquad (1)$$

where $d_i$ is the image-specific depth. This hierarchical encoding stabilizes Elo initialization and reduces the need for downstream human comparisons.

### 2.1.2 CLIP-Based Classification.
For each level in the hierarchy, we compute text and image embeddings using CLIP [13], and measure similarity between image $i$ and prompt $j$ as cosine similarity: $s_{i,l,j} = \frac{\mathbf{v}_i \cdot \mathbf{t}_{l,j}}{||\mathbf{v}_i|| \cdot ||\mathbf{t}_{l,j}||}$. This approach leverages CLIP's zero-shot classification capabilities [11] and builds upon knowledge-enhanced visual models [10].

Classification decisions and confidence scores are derived from these similarities: $c_{i,l} = \arg\max_j s_{i,l,j}$ and $\text{conf}_{i,l} = \frac{\exp(s_{i,l,c_{i,l}}/\tau)}{\sum_j \exp(s_{i,l,j}/\tau)}$ where $\tau = 0.1$ is a temperature parameter controlling the softness of the distribution. Figure 2 shows an example of this process at Level 1 using binary prompts for age grouping. Having established the pre-ordering and initial bucket assignments, we next describe how pairs are selected for human annotation based on uncertainty.

### 2.1.3 Bucket-Aware Rating Initialization.
To obtain a coarse ordering for Elo initialization, the fine-grained groups produced by hierarchical classification are merged into $k$ primary buckets using the mapping $M\colon \{0, 1, \ldots, 2^d - 1\} \to \{0, 1, \ldots, k - 1\}$ defined by $M(g) = \lfloor g\,k/2^d \rfloor$. This rule uniformly distributes groups while preserving their ordinal relationships. We empirically find the optimal number of primary buckets to be $K$ between 3 and 5, which balances ranking accuracy and annotation cost; the smaller $k$ merges overly dissimilar items, while larger $k$ increases unertain cross-bucket comparisons without accuracy benefits.

Accordingly, we set $k = 5$ for CLIP-friendly domains such as FGNET but $k = 3$ for more challenging domains like DHCI and EyePACS, thereby containing noise and comparison overhead. Each image $i$ then receives an Elo rating

$$r_i = r_{\text{base}}(b_i) + \eta_i (1.5 - \text{conf}_i),$$

where $b_i = M(g_i)$ is its bucket, $\eta_i \sim \mathcal{U}(-\delta_b, \delta_b)$ adds controlled randomness, and the confidence term keeps high-confidence samples stable while allowing low-confidence ones to move more freely.

### 2.1.4 Information-Gain-Based Uncertainty and Comparison Prioritization.
We assess the informativeness of comparisons using the KL-divergence from a uniform prior. For items $i$ and $j$ with Elo scores $r_i$ and $r_j$, the pre-comparison distribution is $P_{\text{before}} = [p_{ij}, 1 - p_{ij}]$, where $p_{ij} = \frac{1}{1+10^{(r_j-r_i)/400}}$, following the default setting.

The information gain is computed as:

$$\text{InfoGain}(i, j) = \text{KL}(P_{\text{before}} \| P_{\text{uniform}}) = \sum_k p_k \log \frac{p_k}{0.5} \qquad (2)$$

To prioritize comparisons, we define:

$$\text{Priority}(i, j) = \text{InfoGain}(i, j) \cdot \gamma(b_i, b_j) \cdot \phi(\text{conf}_i, \text{conf}_j) \quad (3)$$

where $\gamma(b_i, b_j) = 1.2$ for cross-bucket pairs (to account for CLIP uncertainty), and 1.0 otherwise; $\phi(\cdot)$ penalizes low-confidence predictions: $\phi(\text{conf}_i, \text{conf}_j) = 2.0 - \text{avg\_conf}$. The uncertainty measure is its complement, normalized by the maximum binary InfoGain (log 2): $\text{uncertainty}(i, j) = 1 - \frac{\text{Priority}(i,j)}{\log 2}$.

### 2.1.5 MergeSort with Uncertainty-Aware Comparison Selection.
Our algorithm follows the *exact* comparison schedule of classical Merge-Sort; the only difference is how each comparison is resolved. Let $\text{uncertainty}(i, j)$ be the KL-based score for items $(i, j)$ and $\theta_t$ an adaptive threshold (Sec. 2.1.5). During each merge, we apply the rule

$$\text{query\_human}(i, j) \iff \text{uncertainty}(i, j) \geq \theta_t, \qquad (4)$$

breaking ties in favor of the human query. If the condition is false, the outcome is decided automatically by $\text{sign}(r_i - r_j)$, where $r$ denotes the current Elo scores.

Because (i) every pair that standard MergeSort examines is still compared, and (ii) deciding (4) takes constant time ($O(1)$), the overall complexity remains $O(n \log n)$. Thus, EZ-Sort preserves the algorithmic optimality while redirecting human effort only to the most uncertain comparisons.

*Adaptive threshold.* The threshold is adapted based on the remaining budget and current accuracy:

$$\theta_t = \theta_0 \Big( 1 + \alpha \frac{\text{remaining}}{\text{total}} \Big) \beta^{\text{accuracy}_t}, \qquad (5)$$

where $t$ denotes the current evaluation cycle (0-based), updated after every batch of human comparisons or every $k$ merge operation, whichever occurs first. Here, $\alpha$ controls budget sensitivity, and $0 < \beta < 1$ encourages increased automation as accuracy improves.

## 3 Experiments and Results

For proof of concept, we evaluated EZ-Sort on three public datasets: FGNET for face age estimation, DHCI for historical image chronology, and EyePACS for retinal image quality. We conducted two types of experiments: (1) inter-rater reliability with expert annotations and (2) annotation efficiency benchmarking across dataset sizes that are common in expert-only scenarios (up to $n = 100$). The reported improvements are statistically significant at $p < 0.05$.

**Inter-rater reliability.** We randomly selected 30 images per dataset and had three domain experts annotate them using absolute classification, sort comparison [4], and EZ-Sort. Results in Table 1 show the inter-rater consistency achieved by each method. Face age estimation (FGNET) achieved uniformly high reliability across all approaches (ICC $\geq$ 0.97), indicating clear visual criteria. DHCI exhibited moderate consistency (ICC = 0.68–0.78), with pairwise methods outperforming classification. For retinal image quality (EyePACS), EZ-Sort attained the highest reliability (ICC = 0.94, Spearman = 0.85), demonstrating robustness in ambiguous medical images.

**Ablation study.** To validate whether hierarchical prompting provides benefit, we analyzed the correlation between the sorted output and the ground truth continuous label, age, which is only

**Table 1: Inter-rater reliability comparison across datasets. Reported metrics: *Sp* (Spearman), *Ke* (Kendall), *Pe* (Pearson), and *ICC* (intraclass correlation). Variance estimates are reported in parentheses where available.**

| Method | Retina (EyePACS) | | | | Historical (DHCI) | | | | Face (FGNET) | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Sp | Ke | Pe | ICC | Sp | Ke | Pe | ICC | Sp | Ke | Pe | ICC |
| Classification | 0.53 (±0.06) | 0.46 (±0.07) | 0.54 (±0.07) | 0.75 (N/A) | 0.39 (±0.06) | 0.33 (±0.05) | 0.42 (±0.06) | 0.68 (N/A) | 0.92 (±0.04) | 0.85 (±0.09) | 0.94 (±0.03) | 0.97 (N/A) |
| Sort comparison [4] | 0.72 (±0.07) | 0.56 (±0.06) | 0.72 (±0.07) | 0.89 (N/A) | 0.47 (±0.17) | 0.35 (±0.15) | 0.47 (±0.17) | **0.78** (N/A) | **0.97** (±0.01) | 0.88 (±0.01) | **0.97** (±0.01) | **0.99** (N/A) |
| EZ-Sort (CIKM) | **0.85** (±0.09) | **0.76** (±0.14) | **0.85** (±0.09) | **0.94** (N/A) | 0.47 (±0.17) | **0.39** (±0.15) | 0.47 (±0.16) | 0.73 (N/A) | 0.96 (±0.01) | **0.91** (±0.02) | 0.96 (±0.01) | 0.99 (N/A) |

**Table 2: The number of human annotations (comparisons) required. FGNET dataset is used.**

| Dataset size | Exhaustive comparison [22] | Sort comparison [4] | EZ-Sort (CIKM) |
|---|---|---|---|
| $n$=30 | 435 | 126 | 89 |
| $n$=50 | 1,225 | 240 | 142 |
| $n$=100 | 4,950 | 582 | 467 |

**Table 3: Dataset characteristics and experimental setup.**

| Dataset | Size | Task | Labels |
|---|---|---|---|
| FGNET | 1,002 | Face Age | 0–69 years (continuous) |
| DHCI | 450 | Historical Dating | 1930s–1970s (5 classes) |
| EyePACS | 28,792 | Retinal Quality | 3-level grading |

available in FGNET. Spearman correlation was 0.90 with EZ-Sort, while flat prompts (with seven class-specific prompts) showed a correlation of 0.83, indicating that hierarchical prompting yields an improvement of 8.4% through progressive refinement.

**Annotation efficiency.** Table 2 shows that EZ-Sort required only 20.5%, 11.6%, and 9.4% of exhaustive comparisons at $n = 30$, 50, and 100, respectively. Compared to sort comparison [4], this corresponds to a relative reduction in human annotation cost of 29.4%, 40.8%, and 19.8% at each respective scale. The most significant gain was observed at $n = 50$, suggesting that EZ-Sort benefits from CLIP pre-ordering most effectively in mid-scale annotation settings. While a slight drop at $n = 100$ reflects increased CLIP uncertainty, our method still maintains substantial efficiency gains over both baselines. These sample sizes reflect real-world-like scenarios commonly found in medical or historical domains. Larger-scale generalization is discussed in Section 4.

**Comparison method allocation.** Human annotation was requested for 23.1%, 18.4%, and 31.2% of comparisons at $n = 30$, 50, and 100, respectively; the rest were resolved automatically using Elo predictions. This demonstrates the adaptive allocation of annotation effort while preserving the $O(n \log n)$ MergeSort structure [36].

**Implementation details.** We used CLIP ViT-B/32 with temperature $\tau$=0.1. Elo ratings $K$=32, $r_{base}$ linearly distributed in [1200, 1800] across $k$ buckets, noise $\delta_b$=75. Adaptive threshold: $\theta_0$=0.15, $\alpha$=0.3, $\beta$=0.9. Priority: $\gamma$=1.2, $\phi$=2.0−avg_conf. Buckets: $k$=5 (FGNET) and $k$=3 (others). Parameters were selected via cross-validation and prompts were generated with GPT-4. CLIP preprocessing requires

39 ms per image on a CPU on average, providing efficient automated pre-ordering prior to human annotation.

**Human annotation cost.** EZ-Sort keeps the theoretical $O(n \log n)$ bound: zero-shot pre-ordering is $(O(kn))$ (with constant $(k)$), bucket-aware Elo is $(O(n))$, and the uncertainty-aware MergeSort follows the canonical $(O(n \log n))$ schedule. Measured against the information-theoretic minimum $(n \ln n)$, our method uses 0.87, 0.73, and 1.01 × that bound for $(n = 30, 50, 100)$; at $(n = 100)$, we require 467 queries versus the 520-query lower limit ($\approx$90% of optimal).

## 4 Discussion

EZ-Sort offers two primary advantages: a significant reduction in human annotation (up to 90.5%) and consistently strong inter-rater reliability across diverse domains. These benefits arise from its hybrid architecture, which combines CLIP-driven priors with uncertainty-aware comparison selection.

Our hierarchical prompting strategy outperforms flat prompting by 2.0 MAE on average, primarily due to CLIP's strength in making binary decisions and facilitating progressive refinement. The KL-based uncertainty prioritizes ambiguous cases, effectively allocating annotation effort where model confidence is lowest.

EZ-Sort has several limitations. First, its performance depends on the reliability of the underlying vision-language model; domain-specific biases in CLIP could affect the initial ranking. Hierarchical prompting, while powerful, may struggle in domains with subtle or poorly defined visual distinctions. Finally, although our simulations suggest scalability to larger datasets, real-world performance in highly imbalanced or noisy settings remains to be validated.

We plan to integrate with annotator reliability models [16] that down-weight noisy labels, which can be crucial in crowd-sourced annotation. Validating the scalability of our approach on larger datasets is a to-do. Additionally, we plan to investigate few-shot fine-tuning [30] or prompt adaptation for VLM models to reduce uncertainty in less common domains. Preliminary trials with Bayesian Elo variants (e.g., TrueSkill [28]) did not yield improvements; thus, exploring new sorting algorithms tailored for uncertainty-guided annotation remains a potential avenue for future direction.

## 5　Gen-AI Usage Disclosure

No AI tools were used in algorithm development, data collection, analysis, hierarchical prompt design, or manuscript content generation. Claude assisted in code debugging and optimization, and GPT-4 provided grammar and stylistic edits for manuscript writing. All core technical and conceptual contributions are original.

## References

[1] Daniel J. Mankowitz, Andrea Michi, Anton Zhernov, Marco Gelmi, Marco Selvi, Cosmin Paduraru, Edouard Leurent, Shariq Iqbal, Jean-Baptiste Lespiau, Alex Ahern, et al. 2023. Faster sorting algorithms discovered using deep reinforcement learning. *Nature* 618, 7964 (2023), 257–263.

[2] Kaiyang Zhou, Jingkang Yang, Chen Change Loy, and Ziwei Liu. 2022. Learning to prompt for vision-language models. *International Journal of Computer Vision* 130, 9 (2022), 2337–2348.

[3] Jing Li, Rafal Mantiuk, Junle Wang, Suiyi Ling, and Patrick Le Callet. 2018. Hybrid-MST: A hybrid active sampling strategy for pairwise preference aggregation. In *Advances in Neural Information Processing Systems*, Vol. 31.

[4] Ikbeom Jang, Garrison Danley, Ken Chang, and Jayashree Kalpathy-Cramer. 2022. Decreasing annotation burden of pairwise comparisons with human-in-the-loop sorting: Application in medical image artifact rating. *arXiv preprint arXiv:2202.04823* (2022).

[5] Ralph Allan Bradley and Milton E. Terry. 1952. Rank analysis of incomplete block designs: I. The method of paired comparisons. *Biometrika* 39, 3/4 (1952), 324–345.

[6] Joseph L. Hansen. 1978. *An Application of the Elo Rating System to Professional Baseball.* Ph.D. Dissertation. Kalamazoo College.

[7] Robin Swezey, Aditya Grover, Bruno Charron, and Stefano Ermon. 2021. Pi-rank: Scalable learning to rank via differentiable sorting. In *Advances in Neural Information Processing Systems*, Vol. 34, 21644–21654.

[8] Xingjian Bai and Christian Coester. 2023. Sorting with predictions. In *Advances in Neural Information Processing Systems*, Vol. 36, 26563–26584.

[9] Lukas Gienapp, Benno Stein, Matthias Hagen, and Martin Potthast. 2020. Efficient pairwise annotation of argument quality. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, 5772–5781.

[10] Sheng Shen, Chunyuan Li, Xiaowei Hu, Yujia Xie, Jianwei Yang, Pengchuan Zhang, Zhe Gan, Lijuan Wang, Lu Yuan, Ce Liu, et al. 2022. K-lite: Learning transferable visual models with external knowledge. In *Advances in Neural Information Processing Systems*, Vol. 35, 15558–15573.

[11] Qi Qian and Juhua Hu. 2024. Online zero-shot classification with CLIP. In *European Conference on Computer Vision*, 462–477. Springer.

[12] Kevin G. Jamieson and Robert Nowak. 2011. Active ranking using pairwise comparisons. In *Advances in Neural Information Processing Systems*, Vol. 24.

[13] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. 2021. Learning transferable visual models from natural language supervision. In *International Conference on Machine Learning*, 8748–8763. PMLR.

[14] Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. 2022. Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation. In *International Conference on Machine Learning*, 12888–12900. PMLR.

[15] Chao Jia, Yinfei Yang, Ye Xia, Yi-Ting Chen, Zarana Parekh, Hieu Pham, Quoc Le, Yun-Hsuan Sung, Zhen Li, and Tom Duerig. 2021. Scaling up visual and vision-language representation learning with noisy text supervision. In *International Conference on Machine Learning*, 4904–4916. PMLR.

[16] Jacob Whitehill, Ting-fan Wu, Jacob Bergsma, Javier R. Movellan, and Paul Ruvolo. 2009. Whose vote should count more: Optimal integration of labels from labelers of unknown expertise. In *Advances in Neural Information Processing Systems*, Vol. 22.

[17] Vikas C. Raykar, Shipeng Yu, Linda H. Zhao, Gerardo H. Valadez, Charles Florin, Luca Bogoni, and Lauren Moy. 2010. Learning from crowds. *Journal of Machine Learning Research* 11 (2010), 1297–1322.

[18] Kaiyang Zhou, Jingkang Yang, Chen Change Loy, and Ziwei Liu. 2022. Conditional Prompt Learning for Vision-Language Models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 16816–16825.

[19] Huazhu Fu, Boyang Wang, Jianbing Shen, Shanshan Cui, Yanwu Xu, Jiang Liu, and Ling Shao. 2019. Evaluation of retinal image quality assessment networks in different color-spaces. In *Medical Image Computing and Computer Assisted Intervention – MICCAI 2019*, 48–56. Springer.

[20] Frank Palermo, James Hays, and Alexei A. Efros. 2012. Dating historical color images. In *Computer Vision – ECCV 2012*, 499–512. Springer.

[21] A. Lanitis, C.J. Taylor, and T.F. Cootes. 2002. Toward automatic simulation of aging effects on face images. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 24, 4 (2002), 442–455.

[22] Louis L. Thurstone. 1927. The method of paired comparisons for social values. *Journal of Abnormal and Social Psychology* 21, 4 (1927), 384.

[23] Xi Chen, Paul N. Bennett, Kevyn Collins-Thompson, and Eric Horvitz. 2013. Pairwise ranking aggregation in a crowdsourced setting. In *Proceedings of the Sixth ACM International Conference on Web Search and Data Mining*, 193–202.

[24] Lucas Maystre and Matthias Grossglauser. 2017. Just sort it! A simple and effective approach to active preference learning. In *International Conference on Machine Learning*, 2344–2353. PMLR.

[25] Hanqing Zhao and Yuehan Luo. 2018. An $O(N)$ Sorting Algorithm: Machine Learning Sort. *arXiv preprint arXiv:1805.04272* (2018).

[26] Ronald L. Graham, Donald E. Knuth, and Oren Patashnik. 1994. *Concrete Mathematics: A Foundation for Computer Science.* Addison-Wesley, 2nd edition.

[27] Thomas H. Cormen, Charles E. Leiserson, Ronald L. Rivest, and Clifford Stein. 2009. *Introduction to Algorithms.* MIT Press, 3rd edition.

[28] Ralf Herbrich, Tom Minka, and Thore Graepel. 2006. TrueSkill™: A Bayesian skill rating system. In *Advances in Neural Information Processing Systems*, Vol. 19.

[29] Akash Kumar Mohankumar and Mitesh Khapra. 2022. Active Evaluation: Efficient NLG Evaluation with Few Pairwise Comparisons. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 8761–8781. Association for Computational Linguistics.

[30] Haokun Liu, Derek Tam, Mohammed Muqeeth, Jay Mohta, Tenghao Huang, Mohit Bansal, and Colin A. Raffel. 2022. Few-shot parameter-efficient fine-tuning is better and cheaper than in-context learning. In *Advances in Neural Information Processing Systems*, Vol. 35, 1950–1965.

[31] Maytal Saar-Tsechansky and Foster Provost. 2004. Active sampling for class probability estimation and ranking. *Machine Learning* 54 (2004), 153–178.

[32] Devichand Budagam, Ashutosh Kumar, Mahsa Khoshnoodi, Sankalp KJ, Vinija Jain, and Aman Chadha. 2024. Hierarchical Prompting Taxonomy: A Universal Evaluation Framework for Large Language Models Aligned with Human Cognitive Principles. *arXiv preprint arXiv:2406.12644* (2024).

[33] Benjamin Samuel Bloom. 1956. *Taxonomy of Educational Objectives: The Classification of Educational Goals.* Longmans, Green.

[34] Lorin W. Anderson, David R. Krathwohl, Peter W. Airasian, Kathleen A. Cruikshank, Richard E. Mayer, Paul R. Pintrich, James Raths, and Merlin C. Wittrock. 2014. *A Taxonomy for Learning, Teaching, and Assessing: A Revision of Bloom's Taxonomy of Educational Objectives.* Pearson.

[35] Grady Booch, Francesco Fabiano, Lior Horesh, Kiran Kate, Jonathan Lenchner, Andrea Loreggia, Keerthiram Murgesan, Nicholas Mattei, Francesca Rossi, et al. 2021. Thinking fast and slow in AI. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 35, 15042–15046.

[36] Richard Cole. 1988. Parallel merge sort. *SIAM Journal on Computing* 17, 4 (1988), 770–785.

[37] Jayashree Kalpathy-Cramer, J. Peter Campbell, Deniz Erdogmus, Peng Tian, Dharanish Kedarisetti, Chace Moleta, James D. Reynolds, Kelly Hutcheson, Michael J. Shapiro, Michael X. Repka, et al. 2016. Plus disease in retinopathy of prematurity: improving diagnosis by ranking disease severity and using quantitative image analysis. *Ophthalmology* 123, 11 (2016), 2345–2351.