

ConSel: Concept-Aware Self-supervised Learning for Regression Beyond Ordinal Tasks

Anonymous CVPR submission

Abstract

001 *Regression is a fundamental problem in computer vision,*
002 *underpinning tasks such as gaze estimation, head pose pre-*
003 *diction, age assessment, aesthetic quality evaluation, crowd*
004 *counting and historical image dating. We introduce **Con-***
005 ***Sel** (Concept-Aware Self-Supervised Regression), a unified*
006 *framework that learns to predict continuous values by pro-*
007 *gressing from coarse semantic concepts to fine-grained nu-*
008 *meric precision. ConSel follows a two-stage curriculum:*
009 *(1) concept-aware self-supervised pretraining, which aligns*
010 *visual embeddings with conceptual guidance through vari-*
011 *ance–covariance regularization without access to ground-*
012 *truth labels, and (2) fine-tuning for precise continuous pre-*
013 *diction. Unlike prior approaches that are optimized only for*
014 *1D ordinal regression, ConSel generalizes to both ordinal*
015 *and multi-dimensional continuous tasks. Evaluated on 15*
016 *benchmark datasets spanning 6 domains, ConSel surpasses*
017 *both domain-specialized and ordinal methods by 15–35%*
018 *while using only 25% of labeled data (4× less than prior*
019 *methods).*

020 1. Introduction

021 In computer vision, regression aims at predicting con-
022 tinuous numerical values directly from images, serving
023 as a foundation for diverse computer-vision tasks such
024 as gaze [6, 8, 89] and head-pose [84] estimation for
025 human–computer interaction, bone-age [52] and disease-
026 severity [81] assessment in healthcare, aesthetic-quality
027 prediction [29] for content scoring in digital media streams,
028 crowd-density estimation [54, 69] in video surveillance, and
029 historical-image dating [49] in historical and archival col-
030 lections. Regression approaches in computer vision can be
031 broadly categorized into two types. (1) **Direct regression**
032 methods and (2) **Vision–language regression** methods.

033 *Direct regression* optimizes convolutional/transformer
034 encoders with standard losses such as Mean-Squared Error
035 (MSE) or smooth-L1. These methods [17, 47, 56, 58, 78,
036 85], neglect the intermediate semantic representations un-
037 derlying continuous variation. While conceptually simple

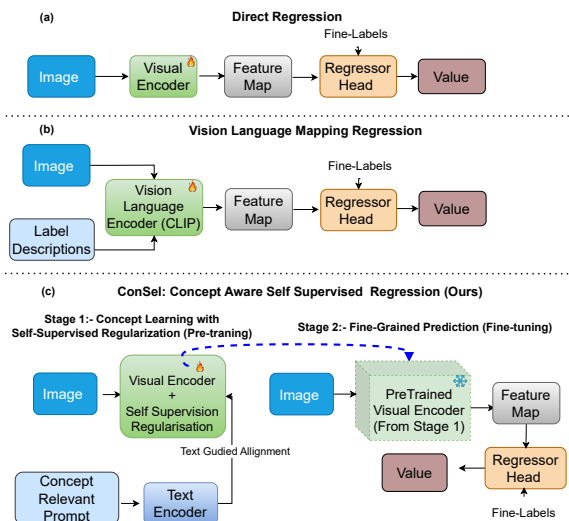


Figure 1. Comparison between (a) **Direct regression**, (b) **VLM-based regression** and (c) our proposed (**ConSel**) regression paradigm.

and widely used, these methods are prone to overfitting with
limited data and generalize poorly to unseen conditions. Re-
cent advances in Vision–Language–Models (VLMs) (e.g.,
CLIP [53], BLIP [34], ALIGN [28]), which have demon-
strated remarkable performance on classification [9, 14] and
segmentation tasks [39, 75], have inspired an alternative
paradigm to regression problem. **Vision–language regres-**
sion methods [13, 37, 73, 79], particularly approaches like
OrdinalCLIP [37] and L2RCLIP [73], convert each label
or rank into a dedicated, learnable prompt (e.g., “age 25
years”, “age 26 years”). This allows VLMs to align visual
features with the semantic order of labels in the text space.
These approaches address the known issues in direct regres-
sion, such as overfitting and poor generalization, especially
where training data is limited. However, converting each
label into a textual description is not scalable, especially
in dense/continuous regression tasks. For example, in tasks
like eye gaze estimation or head pose prediction, where pre-
cision could matter at a scale of 0.5° or less, one would need

057	to generate a very large number of fine grained prompts.	108
058	Building upon this, methods such as L2RCLIP [73] and	109
059	NumCLIP [13] apply the VLMs, particularly CLIP [53], to	110
060	ordinal regression tasks like age estimation and historical	111
061	image dating . These approaches take advantage of the natu-	112
062	ral ordering [21, 71] present in such tasks by grouping labels	113
063	into broader, progressive categories, such as child, young,	114
064	and adult and by using rank or order-aware losses. While,	115
065	these methods excel when the target space exhibits a natural	116
066	one-dimensional ranking (e.g., age groups: child \prec young	117
067	\prec adult), the presupposition of global and natural ordinal-	118
068	ity [21, 51, 71] renders them unsuitable for regression tasks	119
069	where such ranking does not exist, e.g. gaze direction, head	120
070	orientation, or crowd density estimation as shown in Fig 2.	121
071	Empirically, we find that in such settings, ordinal regression	122
072	methods underperform by 15–32% compared to models	123
073	designed specifically for those tasks (see results of Num-	124
074	Clip, l2RCLIP, CLOC and Ordinal-CLIP in Table 1, 2, 5).	125
075	This performance degradation is primarily due to two	126
076	reasons: 1) ordinal losses enforce discrete rank-aware mar-	127
077	gins [51, 73] that are incompatible with manifold or contin-	128
078	uous regression. 2) Additionally, contrastive learning treats	129
079	neighboring labels as negatives, for example, pushing apart	130
080	representations for 24° and 25°. This directly contradicts	131
081	the smoothness prior which is essential to regression [83],	132
082	especially in multi-dimension and continuous tasks. More-	133
083	over, occlusion and real-world variation disrupt the progres-	134
084	sion assumed by group-aware ordinal losses, limiting their	
085	robustness to real world challenges (See Tables 6, 7).	
086	To address the limitations stated above, we propose a	
087	novel strategy to solve regression problems by reformu-	
088	lating them in a concept base curriculum learning [2, 74]	
089	framework ¹ . We first aim to learn high level conceptual	
090	semantics and then progressively refine the model towards	
091	fine-grained target prediction, as illustrated in Figure 1. In-	
092	stead of discretizing labels into ordered bins or enumerating	
093	a prompt for each label, our approach organizes contin-	
094	uous targets into broader concepts, which are refined through	
095	successive curriculum stages. Specifically, our regression	
096	framework follows a two-stage learning process: (a) Con-	
097	cept Learning with Self-Supervised Regularization (pre-	
098	training), and (b) Fine-Grained Prediction (fine-tuning).	
099	In Stage 1 , we pair input images with coarse concep-	
100	tual prompts (e.g., “person looking left,” “person looking	
101	right,” “person looking at center”) that serve as soft regul-	
102	arization anchors during self-supervised pre-training. A cus-	
103	tomized self-supervised loss aligns these soft anchors with	
104	image embeddings while enforcing three representational	
105	constraints: (i) <i>consistency</i> to augmentations of the same	
106	sample, (ii) <i>variance</i> preservation across samples to prevent	
107	collapse, and (iii) <i>covariance</i> minimization across feature	
	dimensions to encourage decorrelation (See Section 3.1 for	108
	details).	109
	In Stage 2 , the visual encoder from stage 1 is frozen,	110
	and a lightweight MLP is trained using only 25% continu-	111
	ous labels (e.g., age values, gaze angle, or crowd density)	112
	refining numeric precision while preserving the learned se-	113
	matic structure. This version of curriculum learning mir-	114
	rors human learning: understanding “left” vs “right” before	115
	predicting precise regression value such as 15.3° yaw and	116
	17.7° pitch.	117
	In a nutshell, our main contributions are as follows:	118
	• We propose a conceptual curriculum learning paradigm	119
	for regression. Our two-stage approach first learns coarse	120
	semantic concepts without exposure to exact numeric la-	121
	bel, then fine-tunes for precise continuous values, im-	122
	proving both performance and generalization.	123
	• We introduce a novel self-supervised pre-training frame-	124
	work for regression that learns semantic structure without	125
	imposing ordinal or rank-aware constraints. A concept-	126
	conditioned variance regularizer links semantic alignment	127
	with statistical consistency, making the approach applica-	128
	ble to both ordinal-style and multidimensional regression.	129
	• We demonstrate strong performance across 6 diverse re-	130
	gression tasks spanning 15 benchmark datasets, surpass-	131
	ing previous methods and achieving consistent gains un-	132
	der distribution shifts and real-world challenges while us-	133
	ing only 25% of the labeled data for fine-tuning.	134
	2. Related work	135
	2.1. Direct and Vision-Language Regression	136
	Traditional regression approaches directly optimize vi-	137
	sual encoders for continuous prediction through end-to-	138
	end learning with standard losses (MSE, MAE). For ordi-	139
	nal regression, methods like CORAL [3], POE [50],	140
	and SORD [12] employ ranking losses to encode order-	141
	ing constraints. Task-specific architectures have been de-	142
	veloped for age estimation [56], gaze prediction [85], and	143
	crowd counting [38]. While effective with abundant la-	144
	beled data, these methods lack semantic structure and strug-	145
	gle with limited supervision or distribution shift. Recent	146
	CLIP-based methods for ordinal regression such as Ordinal-	147
	CLIP [37] and L2RCLIP [73] learn rank-specific prompts;	148
	NumCLIP [13] discretizes targets into 10–20 bins with lin-	149
	guistic descriptions; CLOC [51] employs contrastive rank-	150
	ing losses. While effective, they assume natural 1D order-	151
	ing, making them ill-suited for multi-dimensional tasks like	152
	gaze and head pose estimation where angular relationships	153
	form 2D/3D topological structures [91] rather than linear	154
	progressions as illustrated in Figure 2.	155
	2.2. Curriculum Learning	156
	Curriculum learning improves efficiency by presenting	157
	training examples from easy to hard [2, 74]. Most prior	158

¹Our concept-based curriculum learning is inspired by Definition 3 (Generalized Curriculum Learning) of [74].

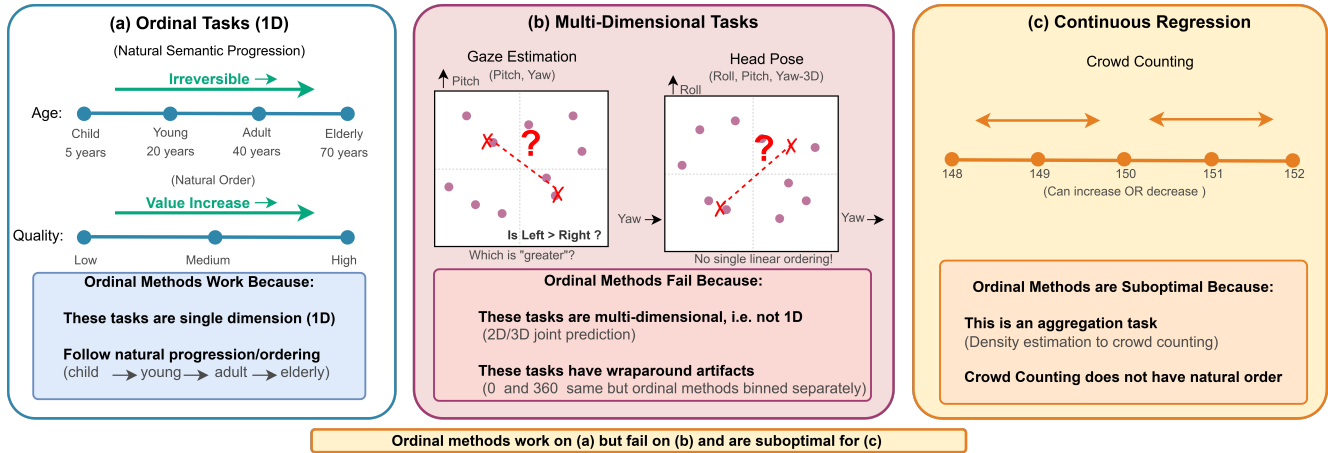


Figure 2. Taxonomy of regression tasks by dimensionality and ordering. Age and image dating are 1D with natural ordinal structure; gaze and head pose are multi-dimensional without inherent ordering. While most Ordinal methods fails on tasks (b) and (c), we show that our proposed method (ConSel) handles all these tasks very well.

work focuses on sample-level ordering [31, 77] or dynamic data selection [19], maintaining fixed loss functions while varying sample selection. We position our method within definition 3 of [74], varying the training criterion across stages, an area still under-explored for regression. Our concept-level curriculum changes the learning objective between stages: Stage-1 employs self-supervised learning with semantic language supervision for concept formation, while Stage-2 uses supervised regression for numeric precision. Unlike NumCLIP’s training with discretized natural bins as target and rank aware contrastive optimisation [13], we use semantic concepts solely for self-supervised structural learning, never as classification targets, avoiding discretization artifacts while establishing concept learning representations.

3. Methodology

Our proposed paradigm teaches regression concepts through a two-stage curriculum that progresses from conceptual understanding to fine-grained prediction. *Stage-1 performs concept-aware self-supervised pretraining* by using our novel *concept-conditioned variance regularization*. These semantics act as soft anchors that organize the space into meaningful regions while preserving intra-concept diversity. *Stage-2 regression fine-tuning* specializes these representations for fine-grained continuous prediction. Figure 1(c) illustrates the overall pipeline of our framework.

3.1. Stage I: Concept Aware Self-Supervised Pre-training

Following the standard practice in non-contrastive self-supervised learning [1, 4, 20, 82], we employ a Siamese architecture with weight-shared encoders as shown in Fig 3. Given an input image x , we generate two augmented views $x_1 = t_1(x)$ and $x_2 = t_2(x)$ using augmentations $t_1, t_2 \sim$

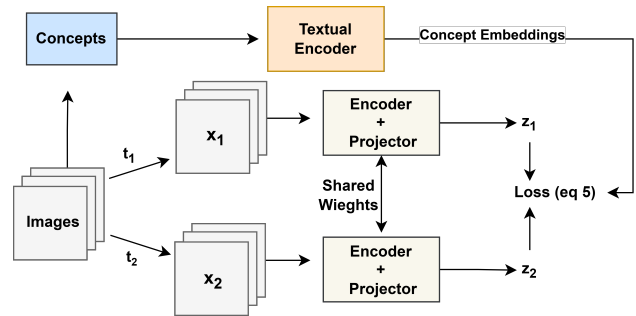


Figure 3. Overview of Stage I pretraining. Note that task relevant *Concepts* can be generated automatically or manually.

\mathcal{T} , where \mathcal{T} includes random crops, color jittering, gaussian blur, and horizontal flips. Both views are processed through an encoder and projector $f_\theta : \mathbb{R}^{H \times W \times 3} \rightarrow \mathbb{R}^d$ to produce embeddings $z_1 = f_\theta(x_1)$ and $z_2 = f_\theta(x_2)$, where d is the embedding dimension. We manually define a minimal, semantically complete concept set $\mathcal{C} = \{c_1, c_2, \dots, c_K\}$ tailored to each regression task. For instance, for directional regression tasks, we define the concept set $\mathcal{C} = \{Looking(left, right, up, down, centre)\}$; for age estimation, we use $\mathcal{C} = \{infant, child, young, adult\}$ (Concept usage and design details are in supplementary Section-B). These concepts are designed to capture coarse-grained semantic variations relevant to the target regression variable, acting as **soft semantic anchors** in the representation space. The concept set \mathcal{C} remains fixed throughout training, providing consistent semantic structure across the dataset.

3.1.1. Loss Formulation

Our pretraining objective integrates conceptual guidance with three complementary regularization terms that prevent representational collapse:

213 **Consistency Loss.** Enforces consistent representations be-
 214 tween augmented views:

$$215 \quad \mathcal{L}_{\text{cons}} = \frac{1}{N} \sum_{i=1}^N \|z_1^i - z_2^i\|_2^2 \quad (1)$$

216 where N is the batch size.

217 **Variance Loss.** To prevent complete collapse where all em-
 218 beddings map to a constant vector, we maintain diversity
 219 across the batch by penalizing low variance for each branch
 220 (z_1, z_2) following [1]:

$$221 \quad \mathcal{L}_{\text{var}} = \frac{1}{d} \sum_{j=1}^d \max(0, \gamma - \sqrt{\text{Var}(Z_j) + \epsilon}) \quad (2)$$

222 where Z_j denotes the j -th embedding dimension across the
 223 batch, $\gamma = 1$ is the target variance threshold, and $\epsilon = 10^{-4}$
 224 ensures numerical stability consistent with [1].

225 **Covariance Loss.** To prevent dimensional collapse where
 226 information concentrates in few dimensions, we decorrelate
 227 embedding dimensions:

$$228 \quad \mathcal{L}_{\text{cov}} = \frac{1}{d} \sum_{i \neq j} [\text{Cov}(Z)]_{ij}^2 \quad (3)$$

229 where $\text{Cov}(Z) \in \mathbb{R}^{d \times d}$ is the covariance matrix of batch
 230 embeddings. This term is particularly important for our cur-
 231 riculum based regression, as it ensures the learned repre-
 232 sentations span a rich, high-dimensional space that supports
 233 fine-grained continuous predictions.

234 **Conceptual Alignment Loss.** We align visual embeddings
 235 with semantic concept embeddings:

$$236 \quad \mathcal{L}_{\text{concept}} = \frac{1}{N} \sum_{i=1}^N \|\bar{z}_i - s_{x_i}\|_2^2 \quad (4)$$

237 where $\bar{z}_i = (z_1^i + z_2^i)/2$ is the mean embedding of both
 238 augmented views. This loss acts as a soft semantic anchor,
 239 structuring the representation space according to human-
 240 interpretable concepts without requiring hard clustering or
 241 negative samples. The complete Stage I objective combines
 242 all terms:

$$243 \quad \mathcal{L}_{\text{Stage I}} = \mathcal{L}_{\text{inv}} + \lambda_v \mathcal{L}_{\text{var}} + \lambda_c \mathcal{L}_{\text{cov}} + \lambda_s \mathcal{L}_{\text{concept}} \quad (5)$$

244 where $\lambda_v = 25$, $\lambda_c = 1.0$, and $\lambda_s = 1.0$ are hyperpa-
 245 rameters balancing the loss components (See Supplemen-
 246 tary Section-B for selection details).

247 All hyperparameters remain fixed throughout training
 248 and are not learned or adapted. The selected values transfer
 249 well across different regression tasks (see Table 1-5).

250 Unlike VICReg [1] which applies variance regu-
 251 larization uniformly, our framework achieves *concept-*
 252 *conditioned variance regularization* through the interplay
 253 of \mathcal{L}_{var} and $\mathcal{L}_{\text{concept}}$: maintaining high intra-concept diver-
 254 sity while ensuring inter-concept separability. This cre-
 255 ates semantically structured embeddings where continuous
 256 regression attributes organize around linguistic concepts,

without requiring contrastive negatives or ordinal ranking
 constraints.

Implicit Ordinal Structure Without Ranking Loss: Our
 method learns ordinal relationships implicitly through the
 interplay of conceptual alignment and variance-covariance
 regularization, without the need for explicit ranking-aware
 losses. First, for 1D ordinal regression (e.g., age, im-
 age aesthetic), semantic concepts from pretrained language
 transformer inherently encode ordering e.g., $\{\textit{infant}, \textit{child},$
 $\textit{young}, \textit{adult}\}$ which have natural sequential relationships
 in text embedding space. By aligning visual embeddings
 to these ordered concept anchors via $\mathcal{L}_{\text{concept}}$, images are
 pulled towards semantically ordered clusters while $\mathcal{L}_{\text{cons}}$ en-
 sures smooth transition between concepts through augmen-
 tation invariance. Specifically, the continuous nature of the
 L2 alignment loss (Eq. 4) allows *representations to interpo-*
late smoothly between concept anchors rather than forming
discrete clusters, supporting fine-grained regression within
 and across ordinal categories. For multi-dimensional re-
 gression (e.g., gaze with pitch-yaw coordinates), the frame-
 work’s effectiveness stems from \mathcal{L}_{cov} and $\mathcal{L}_{\text{concept}}$ ensur-
 ing embedding dimensions remain decorrelated [27]. Con-
 cepts like $\{\textit{left}, \textit{right}, \textit{up}, \textit{down}, \textit{center}\}$ anchor different
 regions of the multi-dimensional embedding space, where
 \mathcal{L}_{var} maintains high within-concept variance to preserve
 fine-grained discriminability along each axis. Thus, our ap-
 proach handles both ordinal 1D and multi-dimensional re-
 gression through semantic guidance and statistical regular-
 ization alone. Figure 4 depicts the learned representation
 spaces for age and eye gaze using UMAP. We report ab-
 lation study results comparing the role and importance of
 different losses with respect to conceptual loss. We also
 perform ablation for evaluating the effect of curriculum in-
 spired paradigm as compared to direct and vision language
 regression (See Tables 8,9).

3.2. Stage II: Regression Fine-tuning

After pretraining converges, we freeze the encoder param-
 eters θ^* and attach a lightweight regression head $h_\psi : \mathbb{R}^d \rightarrow$
 \mathbb{R} . The frozen encoder serves as a powerful feature extrac-
 tor that has learned both semantically-structured and fine-
 grained visual representations during Stage I.

The regression head consists of a 3-layer MLP with hid-
 den dimension 1024, ReLU activation, and dropout ($p =$
 0.2). Given a labeled dataset $\mathcal{D}_{\text{reg}} = \{(x_i, y_i)\}_{i=1}^M$ where
 $y_i \in \mathbb{R}$ are continuous target values, we optimize:

$$257 \quad \mathcal{L}_{\text{Stage II}} = \frac{1}{M} \sum_{i=1}^M (h_\psi(f_{\theta^*}(x_i)) - y_i)^2 \quad (6)$$

This design enables effective transfer learning even when
 $M \ll N$, where N is the pretraining dataset size.

305 4. Experiments

306 We have performed several experiments for evaluating of
307 the overall performance, generalization/transferability and
308 robustness across real world challenges like occlusion and
309 noise. Note that our results reported in this paper were ob-
310 tained using *only 25% of the labeled data during the fine-*
311 *tuning stage*, while the pretraining stage was never exposed
312 to actual labels. Despite using just a quarter of the la-
313 beled data compared to competing methods, our approach
314 achieved stronger performance across all evaluated tasks.

315 4.1. Datasets

316 We used several benchmark datasets; specifically, ETH-
317 XGaze [87], Gaze360 [30], MPIIGaze [86], and EYE-DIAP
318 [18] for eye gaze estimation, while 300W-LP [93] BIWI
319 [16] and AFLW2000 [92] were used for head pose estima-
320 tion. For the latter, we trained our framework on 300W-
321 LP and tested on BIWI and AFLW2000 following previ-
322 ous work [23]. **For crowd counting**, we utilized four
323 benchmark datasets: UCF-QNRF [26], UCF-CC-50 [25],
324 ShanghaiTech Part A [88], and JHU-CROWD++ [64]. We
325 adopted standard train/test splits and reported Mean Abso-
326 lute Error (MAE) as the primary evaluation metric follow-
327 ing [44]. To evaluate generalization of our method, we
328 additionally performed a cross-domain evaluation by train-
329 ing on UCF-QNRF and ShanghaiTech Part A and testing
330 on the remaining datasets. For **1D Ordinal regression**
331 tasks we used the MORPH II [55] and Adience [15] for age
332 estimation, Crowd Beauty[60] for image aesthetic quality
333 prediction and HCI dataset [49] for Historical Image Dat-
334 ing (HID). Following [35, 37, 57, 61] only 5,492 images
335 of Caucasian descent were used from MORPH II to mini-
336 mize cross-race interference. The Adience dataset contains
337 26,580 color images of 2,284 subjects, each assigned to one
338 of eight age groups; we evaluated on it using 5 folds as there
339 is no original train/test split. More details on the statistics
340 and preprocessing of all datasets are provided in the supple-
341 mentary material Section-A.

342 4.2. Implementation Details

343 For all regression tasks, we used the same two-stage
344 training procedure with manually created task-specific
345 concept vocabularies (e.g., “looking left/right” for gaze,
346 “crowded/sparse” for crowd counting, “young/elderly” for
347 age). We randomly initialized ResNet-18 [22] to avoid
348 domain-specific biases from ImageNet pretraining, en-
349 abling task-agnostic concept learning particularly beneficial
350 for tasks requiring spatial reasoning (crowd counting) or
351 fine-grained analysis (age/gaze) rather than object recogni-
352 tion. In the first stage, we employ OREPA [24] for efficient
353 training of ResNet-18 alongside frozen DistilBERT [59],
354 which provides semantic concept representations. We stan-
355 dardized both visual and textual feature dimensions to 1024

Table 1. Comparison with the state-of-the-art methods on eye gaze as well as with ordinal classification methods using domain generalization. The best results are in bold and the second-best are underlined. E:ETH-XGaze, G:Gaze360, M:MPIIGaze and D:EYE-DIAP. Results are reported as angular error (\downarrow lower is better). Results marked with * were obtained using authors’ code.

Methods	Train \rightarrow Test				Avg
	E \rightarrow M	E \rightarrow D	G \rightarrow M	G \rightarrow D	
PureGaze [7]	7.08	7.48	9.28	9.32	8.29
CLIP-Gaze[80]	<u>6.41</u>	7.51	6.89	7.06	6.97
LG-Gaze [79]	6.45	<u>7.22</u>	6.83	<u>6.86</u>	6.84
CDG [76]	6.73	7.95	7.03	7.27	7.25
CoOp [90]	7.44	7.42	7.41	7.15	7.36
CLOC [51]	9.87	11.27	9.93	10.89	10.49
MoCo-v3 [5]	9.43	10.78	9.12	10.01	9.84
DINov2 [48]	8.02	10.14	8.17	9.36	8.92
VICReg [1]	7.93	8.02	7.45	8.54	7.99
Ordinal-CLIP* [37]	8.08	10.56	9.10	9.91	9.41
Num-CLIP* [13]	8.03	10.07	8.05	9.26	8.85
L2RCLIP* [73]	8.51	9.94	8.31	9.08	8.96
Ours	6.19	6.97	6.41	6.82	6.59

for effective alignment and trained for 200 epochs with
batch size 32 and all images were resized to 224 \times 224. Once
concept learning converges, we froze the learned ResNet-
18 visual encoder and trained only a lightweight 3-layer
MLP regressor head for continuous value prediction, isolat-
ing concept grounding from task-specific regression. The
second stage is trained for 50 epochs and with just 25% of
stratified labeled data from training set. Early stopping is
used to ensure optimal generalization.

4.3. Comparison on Multi-Dimensional Regression

For gaze estimation, ETH-XGaze (E) and Gaze360 (G)
served as training datasets, and MPIIGaze (M) and EYE-
DIAP (D) were used for testing to ensure fairness with prior
cross-domain evaluation protocols [7, 79, 80]. Accord-
ingly, results were evaluated across four cross-domain set-
tings (E \rightarrow M, E \rightarrow D, G \rightarrow M, G \rightarrow D). Table 1 reports cross-
domain gaze estimation results using mean angular error in
degrees (MAE; lower is better).

Our method achieved the lowest average error (**6.59 $^\circ$**)
across all four train \rightarrow test settings, outperforming both spe-
cialized gaze models [7, 76, 79, 80] and state-of-the-art
VLM/CL based ordinal regression methods [13, 37, 51,
73, 90]. Table 1 also shows that VLM/CL-based ordi-
nal regression networks underperform on eye-gaze estima-
tion compared to domain-specific methods (e.g., LG-Gaze
model [79]). The same holds true when compared against
our method, where VLM/CL networks Ordinal-CLIP [37],
Num-CLIP [13] and CLOC [51] underperform by a mar-
gin of 42.8%, 34.3%, and 59.2%, respectively. Beyond
cross-dataset evaluation, we compared with PureGaze [7]
and CLIPGaze [80] to assess performance under unsuper-
vised domain adaptation (UDA), where the model is trained

Table 2. Head pose estimation on BIWI and AFLW2000. The best is in **bold** and the second-best is underlined. Results marked with * were obtained using authors’ code.

Method	BIWI (MAE ↓)	AFLW2000 (MAE ↓)
Yang [78]	4.61	5.34
RankPose [10]	4.40	3.74
Valle [68]	<u>3.66</u>	3.83
Hempel [23]	3.78	<u>3.61</u>
Nguyen [46]	9.00	7.23
Musti [45]	6.72	7.82
Tourani [67]	4.47	5.28
VICReg [1]	4.21	4.30
MoCo-v3 [5]	4.62	4.27
DINOv2 [48]	4.41	4.22
CLOC[51]	4.55	4.30
Num-CLIP* [13]	6.05	7.13
L2RCLIP* [73]	7.34	7.89
Ours	3.12	3.27

388 on labeled source data and a small set of unlabeled target samples. Our approach outperforms all UDA methods
 389 (See Supplementary Section-C). For headpose estimation,
 390 following the cross-dataset protocol of [23], we trained on
 391 300W-LP and evaluated on BIWI and AFLW2000 as shown
 392 in Table 2. Our method achieved 3.12° MAE on BIWI

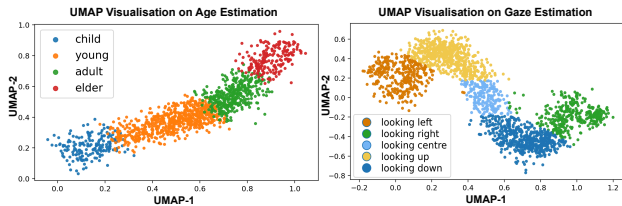


Figure 4. UMAP plot of learned representations. **Left** Age embeddings show smooth progression from young to elderly, validating implicit ordinal structure. **Right:** Gaze embeddings form a wave pattern capturing the continuous topology of directional space guided by semantic concepts. More visualizations in supplementary.

393 and 3.27° MAE on AFLW2000 surpassing other competitors.
 394 Valle [68] obtained 3.66° on BIWI, while Hempel
 395 [23] achieved 3.61° on AFLW2000. Self-supervised approaches including MoCo-v3 [5] and DINOv2 [48] ranged from 4.2–4.6° MAE. Vision-language methods using ordinal representations (NumCLIP [13], L2RCLIP [73]) obtained 6–7.9° MAE.

4.4. Comparison on 1D-Regression Tasks

401 **Age-Estimation:** Table 3 (col: 2, 3) show results on MORPH II and Adience dataset. Our method achieved 1.89 MAE on MORPH II (9% better than NumCLIP). On Adience, we obtained 76.0% Accuracy and 0.32 MAE, exceeding the best prior accuracy (69.4% from NumCLIP), and reducing MAE compared to the strongest prior methods (0.36 from L2RCLIP) by 0.04 (9.8% better). Notably, our method achieved these results despite being designed to handle both single-dimensional (age, aesthetics) and multi-dimensional

Table 3. Combined results on MORPH II, Adience, and HCI. The best is in **bold**; the second-best is underlined. **Blue numbers** = relative (%) improvement over the second-best.

Method	MORPH II	Adience		HCI (HID)	
	MAE ↓	Acc ↑	MAE ↓	Acc ↑	MAE ↓
DRC-ORID [32]	2.26	–	–	–	–
POE [36]	2.35	60.5	0.47	54.68	0.67
PML [11]	2.31	–	–	–	–
MWR [62]	2.13	62.6	0.45	57.80	0.58
CNNPOR [42]	–	57.4	0.55	50.12	0.82
GP-DNNOR [43]	–	57.4	0.54	46.60	0.76
SORD [12]	–	59.6	0.49	53.40	0.70
GOL [33]	–	62.5	0.43	56.20	0.55
RnC [83]	–	46.8	0.86	56.64	0.72
CoOp [90]	2.39	60.6	0.50	51.90	0.76
VICReg [1]	2.26	<u>68.3</u>	0.46	68.53	0.38
MoCo-v3 [5]	2.29	63.2	0.51	57.53	0.72
DINOv2 [48]	2.26	65.9	0.49	62.31	0.66
OrdinalCLIP [37]	2.32	61.2	0.47	56.44	0.67
CLOC [51]	2.27	63.0	0.41	62.08	0.55
L2RCLIP [73]	2.11	66.2	<u>0.36</u>	67.22	0.43
NumCLIP [13]	<u>2.08</u>	–	–	<u>69.61</u>	<u>0.35</u>
Ours	1.89	76.0	0.32	78.51	0.26
	+9.13%	+11.27%	+11.11%	+12.79%	+25.71%

(gaze, head pose) tasks with the same framework, whereas prior methods demonstrated their performance only on 1D ordered outputs.

Image Quality Estimation: For image quality estimation, we evaluated on two different ordinal tasks 1) image dating (HCI dataset) and 2) aesthetics scoring (CrowdBeauty). These tasks also benefited from our semantic concept guidance, as concepts naturally encode ordering (e.g., *old* vs. *recent* for dating). Our method achieved 78.51% accuracy and 0.26 MAE on HCI (See Table 3; col 4), significantly outperforming NumCLIP (69.61% accuracy, 0.35 MAE). This represents an 8.90% point accuracy gain and 25.7% MAE reduction, demonstrating that semantic concepts (*ancient*, *vintage*, *modern*, *recent*) provide ordinal structure without explicit ranking loss. On CrowdBeauty (Table 4), we achieved 79.62% overall accuracy and 0.209 MAE, improving over NumCLIP (76.53%, 0.236) by +3.09 points and 11.4% MAE reduction. Category-wise, our method excelled on Nature (83.12%), Urban (81.47%), and Animal (76.39%) categories, validating that semantic guidance transfers across diverse visual domains. Performance on People (77.52%) remains competitive, though slightly behind specialized method. Recently, Wang et al. [72] reported results on ordinal tasks (HCI, Crowd Beauty, and Adience). Despite using 100% of the training data (4 times ours), their performance on HCI and Crowd Beauty remains below ours, though they achieve a superior 0.12 MAE on the Adience age-group classification task.

4.5. Comparison on Crowd Estimation

For crowd estimation, we compared against specialized crowd density estimators (P2R [41], P2PNet [65], DM-

Table 4. Category-wise results on Image Aesthetic Quality Task. Accuracy (%) \uparrow and MAE \downarrow .

Method	Accuracy (%) \uparrow					MAE \downarrow				
	Nature	Animal	Urban	People	Overall	Nature	Animal	Urban	People	Overall
CNNPOR [42]	71.86	69.32	69.09	69.94	70.05	0.294	0.322	0.325	0.321	0.316
SORD [12]	73.59	70.29	73.25	70.59	72.03	0.271	0.308	0.276	0.309	0.290
POE [36]	73.62	71.14	72.78	72.22	72.44	0.273	0.299	0.281	0.293	0.287
GOL [33]	73.80	72.40	74.20	69.60	72.70	0.270	0.280	0.260	0.310	0.280
CoOp [90]	72.74	71.46	72.14	69.34	71.42	0.285	0.298	0.294	0.330	0.302
OrdinalCLIP [37]	73.65	72.85	73.20	72.50	73.05	0.273	0.279	0.277	0.291	0.280
Ord2Seq [37]	<u>78.09</u>	<u>75.74</u>	72.83	69.24	74.43	<u>0.225</u>	0.257	0.275	0.319	0.264
L2RCLIP [73]	73.51	75.26	77.76	78.69	76.07	0.267	0.253	0.216	0.246	0.245
NumCLIP [13]	75.20	75.24	<u>79.49</u>	76.17	<u>76.53</u>	0.249	<u>0.250</u>	<u>0.208</u>	0.238	<u>0.236</u>
Ours	83.12	76.39	81.47	<u>77.52</u>	79.62	0.187	0.213	0.199	<u>0.239</u>	0.209

Count [70]). Table 5 shows that on four dataset, ShanghaiTech Part A (A), UCF-QNRF (B), UCF-CC-50 (C), and JHU-CROWD++ (Q), our method attained the lowest MAE (lower is better): 46.8 (A), 65.8 (B), 200.1 (C), and 51.4 (Q) respectively. While purely ordinal regression (OrdinalCLIP) methods perform poorly across splits, bin-based adaptations (Num-CLIP, CoOp) close some of the gap but still lag behind specialized methods (e.g., P2R, MPCount) and also lag behind our method. Moreover, we also performed cross-domain evaluation surpassing the prior methods (see supplementary section-C).

Table 5. Crowd counting results (MAE; \downarrow lower is better) using standard splits: ShanghaiTech Part A (A), UCF-QNRF (B), UCF-CC-50 (C), and JHU-CROWD++ (Q). Best is in **bold**; second best is underlined.

Methods	A(Shang)	B(UCF)	C(50)	Q(JHU)
BL [44]	62.8	88.7	229.3	-
P2R [41]	<u>51.02</u>	83.2	-	<u>58.83</u>
P2PNet [65]	52.74	85.3	-	61.25
SASNet [66]	67.0	-	258.4	-
ChIL [63]	57.5	80.3	-	57.0
MAN [40]	56.8	77.3	-	53.4
DMCount [70]	59.7	85.6	<u>211.0</u>	-
Ordinal-CLIP [37]	59.9	82.8	244.4	60.5
Num-CLIP [13]	54.2	<u>73.7</u>	221.3	61.3
CLOC [51]	67.9	94.8	255.9	67.2
CoOp [90]	64.0	75.2	219.1	62.5
Ours	46.8	65.8	200.1	51.4

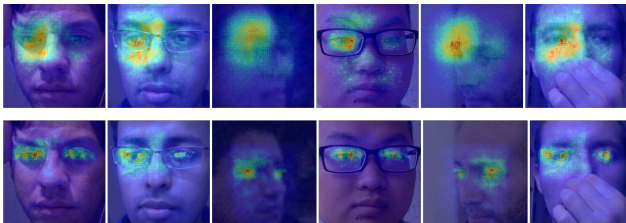


Figure 5. Grad-CAM++ visualisation on challenging samples of eye-gaze. Row 1: after Stage 1 (concept pretraining). Row 2: after Stage 2 (fine-tuning on precise labels).

4.6. Robustness to Distribution Shift

Cross-domain evaluation on eye gaze and crowd counting validates the generalizability of our method across di-

Table 6. MAE \downarrow under distribution shift on MORPH II. Best is **bold**, second-best is underlined.

Method	number of reduced classes - % age of reduced samples/ class							
	10-80	10-90	20-80	20-90	30-80	30-90	40-80	40-90
CoOp [90]	2.71	2.85	2.98	3.51	3.06	3.36	2.99	3.30
OrdinalCLIP [37]	2.61	2.67	2.77	3.06	2.86	3.21	2.84	3.12
L2RCLIP	2.28	2.30	2.37	2.43	2.51	2.61	2.68	2.79
(Ours)	2.03	2.12	2.14	2.19	2.24	2.28	2.33	2.41

verse regression tasks. To further assess robustness under distribution shift, we conducted experiments on MORPHII for age estimation. Following the protocol of [37], we randomly selected several age ranges (labels) and removed a large portion of samples from those ranges, creating an imbalanced training distribution that differs significantly from the test distribution. Table 6 shows our method consistently outperformed baselines across all distribution shift scenarios. Under severe imbalance (40 classes reduced by 90%), our approach achieved 2.41 MAE compared to L2RCLIP's 2.79 and OrdinalCLIP's 3.12, representing **13.6%** and **22.8%** improvement respectively. Even under extreme shifts (10 classes, 90% reduction), we maintained 2.12 MAE versus 2.30 for L2RCLIP. Our robustness emerges from two synergistic factors: semantic anchors encode generalizable structure agnostic to training distribution, while covariance regularization maintains rich, high-dimensional representations that generalize effectively to underrepresented ranges.

4.7. Robustness to Real-World Challenges

We evaluated robustness on natural variations: gaze estimation with accessories (glasses/caps) excluded from training, and age estimation with 10% facial occlusion. Table 7 demonstrates substantial advantages: 5.80° gaze error versus NumCLIP's 13.76° (57.9% improvement), and 4.45/3.38 years MAE on occluded age tasks versus 6.01/5.81 (26-42% better). Figure 5 depicts GradCAM++ attention on challenging samples with hand and glasses occlusions, showing that our two-stage curriculum enables the model to focus on discriminative regions despite these real world challenges.

Table 7. Robustness results: for gaze estimation, we evaluated on glasses/cap subset (angular MAE in degrees) after fine-tuning without those samples. For age estimation, we applied a 10% central-face occlusion on MORPH II and Adience. Results marked with * are obtained using authors’ code.

Method	Gaze (deg) ↓	MORPH II (y) ↓	Adience ↓
Ours	5.80	4.45	3.38
CLOC [51]	16.98	6.32	5.33
NumCLIP* [13]	13.76	6.01	5.81
L2RCLIP* [73]	14.01	6.77	5.30
Ordinal-CLIP* [13]	14.54	7.33	6.11

Table 8. Ablation on loss components in the curriculum paradigm. Eye Gaze (G⇒M) is mean angular error (↓°); Age (MORPH II) and HID are MAE (↓). **Blue Numbers** are the relative improvement vs the row above. Results are mean±std over 3 seeds.

Configuration	Eye Gaze (↓°)	Age (↓)	HID (↓)
Consistency loss only	10.24±1.97	4.21±0.73	0.69±0.35
	(—)	(—)	(—)
+ Concept loss	7.63±1.18	2.30±0.41	0.42±0.27
	(+25.5%)	(+45.4%)	(+39.1%)
+ Var/Cov (no concept)	7.45±0.35	2.26±0.23	0.38±0.14
	(+2.4%)	(+1.7%)	(+9.5%)
Full	6.19±0.27	1.89±0.15	0.26±0.06
	(+16.9%)	(+16.4%)	(+31.6%)

4.8. Ablation Studies

We conducted ablation studies to validate our design choices, analyzing both loss components (Table 8) and training paradigms (Table 9). We evaluated on three diverse tasks: multi-dimensional gaze regression, 1D ordinal age regression (MORPH II), and on HCI. We systematically added loss terms to understand their individual and synergistic effect within the curriculum learning framework.

Using only $\mathcal{L}_{\text{cons}}$ yielded the weakest performance (10.24° gaze error, 4.21 years age MAE, 0.69 HID MAE), as the model lacked both semantic structure and collapse prevention mechanisms. Without variance-covariance regularization, representations may suffer from dimensional collapse; without conceptual guidance, the learned space has no semantic organization.

Introducing $\mathcal{L}_{\text{concept}}$ produced substantial gains: **25.5%** improvement on gaze, **45.4%** on age, and **39.1%** on HID. The particularly large gain on age estimation validates our core claim that semantic concepts (*infant*→*adult*) provided implicit ordinal structure without explicit ranking losses. For HID image dating, temporal concepts (*ancient*→*modern*) similarly encode natural ordering.

Replacing conceptual alignment with \mathcal{L}_{var} and \mathcal{L}_{cov} (Table 8 row 3) yielded comparable performance to concept-only (row 2) on age and HID, but slightly better on gaze (7.45° vs. 7.63°). This suggests that for multi-dimensional regression, decorrelating embedding dimensions is crucial for encoding independent axes (pitch vs. yaw), while for 1D ordinal tasks, semantic structure dominates.

Full Model (All Losses): Combining all four terms achieved the best performance across all tasks, with ad-

Table 9. Ablation across training regimes. Eye Gaze (G⇒M) is reported as mean angular error (↓°); MORPH II as MAE (↓ years). Δ_{prev} is the relative improvement vs the row above. Results are mean±std over 3 seeds.

Training	Eye Gaze (↓Err, °)		MORPHII (↓ MAE)	
	MAE Err	$\Delta_{\text{prev}}\%$	MAE	$\Delta_{\text{prev}}\%$
Direct	10.07 ± 2.58	—	5.87 ± 1.04	—
+Text (E2E)	8.04 ± 1.39	+20.16%	3.89 ± 0.95	+33.73%
+Cur (Ours)	6.19 ± 0.27	+23.01%	1.89 ± 0.15	+51.41%

Definitions: *Direct (Single Stage)* = no curriculum, no text guidance. *+Text (E2E)* = end-to-end with textual concepts and task labels. *+Cur (Ours)* = two-stage with concept guidance in Stage-1.

ditional gains of **16.9%** (gaze), **16.4%** (age), and **31.6%** (HID) over the previous configuration. This demonstrates that conceptual alignment and variance-covariance regularization are *complementary* rather than redundant: concepts provide semantic structure while var-cov prevents collapse and maintains high-dimensional expressivity.

Training Paradigm Analysis (Table 9): Progressive enhancements demonstrate clear cumulative benefits: **Direct regression** (no text) produced scattered embeddings with no structure; **+Text guidance** improved results by 20-34%; **+Curriculum** (ours) gained additional 23-51% with the most structured embeddings; eye gaze forming clear progressive clusters and age showing smoothest ordinal progression (Fig. 4). Furthermore, the significantly lower standard deviation ($\pm 0.27^\circ$ vs. $\pm 2.58^\circ$) indicates enhanced training stability. We also conducted experiments with irrelevant concepts (semantically valid but task-unrelated) and LLM-generated concepts (automatically produced via GPT-4). Results confirmed that LLM-generated concepts also achieved good performance demonstrating automation potential (See Supplementary Section-C).

5. Conclusion

We introduced ConSel, a new framework designed to handle ordinal and non-ordinal regression through language-guided regularization. ConSel first learned conceptual representations via non-contrastive self-supervised pretraining with language alignment, allowing it to build meaningful ordinal structures from semantic language features without the need for explicit ranking losses. This approach enabled ConSel to achieve strong performance on both ordinal and non-ordinal tasks. The effectiveness of ConSel is dependent on the choice of concepts. Strong results are obtained when the concepts are semantically relevant for the task. Notably, ConSel demonstrated significant advantage in low-data regimes, outperforming its competitors while using 4x less training data. This highlights ConSel’s potential for efficient learning in scenarios with limited data availability.

557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593
594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613

References

- [1] Adrien Bardes, Jean Ponce, and Yann LeCun. Vircreg: Variance-invariance-covariance regularization for self-supervised learning. In *International Conference on Learning Representations*, 2022. 3, 4, 5, 6
- [2] Yoshua Bengio, Jérôme Louradour, Ronan Collobert, and Jason Weston. Curriculum learning. In *Proceedings of the 26th Annual International Conference on Machine Learning*, pages 41–48, 2009. 2
- [3] Wenzhi Cao, Vahid Mirjalili, and Sebastian Raschka. Rank consistent ordinal regression for neural networks with application to age estimation. *Pattern Recognition Letters*, 140: 325–331, 2020. 2
- [4] Xinlei Chen and Kaiming He. Exploring simple siamese representation learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 15750–15758, 2021. 3
- [5] Xinlei Chen, Saining Xie, and Kaiming He. An empirical study of training self-supervised vision transformers. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 9640–9649, 2021. 5, 6
- [6] Yihua Cheng and Feng Lu. Dvgaze: Dual-view gaze estimation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 20632–20641, 2023. 1
- [7] Yihua Cheng, Yiwei Bao, and Feng Lu. Puregaze: Purifying gaze feature for generalizable gaze estimation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 436–443, 2022. 5
- [8] Yihua Cheng, Hengfei Wang, Zhongqun Zhang, Yang Yue, Boeun Kim, Feng Lu, and Hyung Jin Chang. 3d prior is all you need: Cross-task few-shot 2d gaze estimation. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 23891–23900, 2025. 1
- [9] Alessandro Conti, Enrico Fini, Massimiliano Mancini, Paolo Rota, Yiming Wang, and Elisa Ricci. Vocabulary-free image classification. *Advances in Neural Information Processing Systems*, 36:30662–30680, 2023. 1
- [10] Donggen Dai, Wangkit Wong, and Zhuojun Chen. Rankpose: Learning generalised feature with rank supervision for head pose estimation. *arXiv preprint arXiv:2005.10984*, 2020. 6
- [11] Zongyong Deng, Hao Liu, Yaoxing Wang, Chenyang Wang, Zekuan Yu, and Xuehong Sun. Pml: Progressive margin loss for long-tailed age classification. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10503–10512, 2021. 6
- [12] Raul Diaz and Amit Marathe. Soft labels for ordinal regression. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4738–4747, 2019. 2, 6, 7
- [13] Yao Du, Qiang Zhai, Weihang Dai, and Xiaomeng Li. Teach clip to develop a number sense for ordinal regression. In *European Conference on Computer Vision*, pages 1–17. Springer, 2024. 1, 2, 3, 5, 6, 7, 8
- [14] Lisa Dunlap, Yuhui Zhang, Xiaohan Wang, Ruiqi Zhong, Trevor Darrell, Jacob Steinhardt, Joseph E Gonzalez, and Serena Yeung-Levy. Describing differences in image sets with natural language. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 24199–24208, 2024. 1
- [15] Eran Eidinger, Roei Enbar, and Tal Hassner. Age and gender estimation of unfiltered faces. *IEEE Transactions on information forensics and security*, 9(12):2170–2179, 2014. 5
- [16] Gabriele Fanelli, Matthias Dantone, Juergen Gall, Andrea Fossati, and Luc Van Gool. Random forests for real time 3d face analysis. *International journal of computer vision*, 101(3):437–458, 2013. 5
- [17] Tobias Fischer, Hyung Jin Chang, and Yiannis Demiris. Rtgene: Real-time eye gaze estimation in natural environments. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 334–352, 2018. 1
- [18] Kenneth Alberto Funes Mora, Florent Monay, and Jean-Marc Odobez. Eyediap: A database for the development and evaluation of gaze estimation algorithms from rgb and rgb-d cameras. In *Proceedings of the symposium on eye tracking research and applications*, pages 255–258, 2014. 5
- [19] Alex Graves, Marc G Bellemare, Jacob Menick, Remi Munos, and Koray Kavukcuoglu. Automated curriculum learning for neural networks. In *International Conference on Machine Learning*, pages 1311–1320. PMLR, 2017. 3
- [20] Jean-Bastien Grill, Florian Strub, Florent Altché, Corentin Tallec, Pierre Richemond, Elena Buchatskaya, Carl Doersch, Bernardo Avila Pires, Zhaohan Guo, Mohammad Gheshlaghi Azar, et al. Bootstrap your own latent: A new approach to self-supervised learning. In *Advances in Neural Information Processing Systems*, pages 21271–21284, 2020. 3
- [21] Pedro Antonio Gutiérrez, Maria Perez-Ortiz, Javier Sanchez-Monedero, Francisco Fernandez-Navarro, and Cesar Hervás-Martinez. Ordinal regression methods: survey and experimental study. *IEEE Transactions on Knowledge and Data Engineering*, 28(1):127–146, 2015. 2
- [22] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. 5
- [23] Thorsten Hempel, Ahmed A Abdelrahman, and Ayoub Al-Hamadi. Toward robust and unconstrained full range of rotation head pose estimation. *IEEE Transactions on Image Processing*, 33:2377–2387, 2024. 5, 6
- [24] Mu Hu, Junyi Feng, Jiashen Hua, Baisheng Lai, Jianqiang Huang, Xiaojin Gong, and Xian-Sheng Hua. Online convolutional re-parameterization. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 568–577, 2022. 5
- [25] Haroon Idrees, Imran Saleemi, Cody Seibert, and Mubarak Shah. Multi-source multi-scale counting in extremely dense crowd images. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2547–2554, 2013. 5
- [26] Haroon Idrees, Muhammad Tayyab, Kishan Athrey, Dong Zhang, Somaya Al-Maadeed, Nasir Rajpoot, and Mubarak Shah. Composition loss for counting, density map estimation and localization in dense crowds. In *Proceedings of the European conference on computer vision (ECCV)*, pages 532–546, 2018. 5

614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671

- 672 [27] Saahil Jain, Ashish Jain, and Joost van de Weijer. Bar-
673 low twins for multimodal self-supervised learning. *arXiv*
674 *preprint arXiv:2112.09061*, 2021. 4
- 675 [28] Chao Jia, Yinfei Yang, Ye Xia, Yi-Ting Chen, Zarana Parekh,
676 Hieu Pham, Quoc Le, Yun-Hsuan Sung, Zhen Li, and Tom
677 Duerig. Scaling up visual and vision-language representa-
678 tion learning with noisy text supervision. In *International*
679 *conference on machine learning*, pages 4904–4916. PMLR,
680 2021. 1
- 681 [29] Junjie Ke, Keren Ye, Jiahui Yu, Yonghui Wu, Peyman Mi-
682 lanfar, and Feng Yang. Vila: Learning image aesthetics from
683 user comments with vision-language pretraining. In *Pro-*
684 *ceedings of the IEEE/CVF Conference on Computer Vision*
685 *and Pattern Recognition*, pages 10041–10051, 2023. 1
- 686 [30] Petr Kellnhofer, Adria Recasens, Simon Stent, Wojciech Ma-
687 tusik, and Antonio Torralba. Gaze360: Physically uncon-
688 strained gaze estimation in the wild. In *Proceedings of*
689 *the IEEE/CVF international conference on computer vision*,
690 pages 6912–6921, 2019. 5
- 691 [31] M Pawan Kumar, Benjamin Packer, and Daphne Koller. Self-
692 paced learning for latent variable models. In *Advances in*
693 *Neural Information Processing Systems*, 2010. 3
- 694 [32] Seon-Ho Lee and Chang-Su Kim. Deep repulsive clustering
695 of ordered data based on order-identity decomposition. In *Inter-*
696 *national Conference on Learning Representations*, 2020.
697 6
- 698 [33] Seon-Ho Lee, Nyeong Ho Shin, and Chang-Su Kim. Geo-
699 metric order learning for rank estimation. *Advances in Neu-*
700 *ral Information Processing Systems*, 35:27–39, 2022. 6, 7
- 701 [34] Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi.
702 Blip: Bootstrapping language-image pre-training for unified
703 vision-language understanding and generation. In *Interna-*
704 *tional conference on machine learning*, pages 12888–12900.
705 PMLR, 2022. 1
- 706 [35] Wanhua Li, Jiwen Lu, Jianjiang Feng, Chunjing Xu, Jie
707 Zhou, and Qi Tian. Bridgenet: A continuity-aware proba-
708 bilistic network for age estimation. In *Proceedings of*
709 *the IEEE/CVF Conference on Computer Vision and Pattern*
710 *Recognition*, pages 1145–1154, 2019. 5
- 711 [36] Wanhua Li, Xiaoke Huang, Jiwen Lu, Jianjiang Feng,
712 and Jie Zhou. Learning probabilistic ordinal embed-
713 dings for uncertainty-aware regression. In *Proceedings of*
714 *the IEEE/CVF conference on computer vision and pattern*
715 *recognition*, pages 13896–13905, 2021. 6, 7
- 716 [37] Wanhua Li, Xiaoke Huang, Zheng Zhu, Yansong Tang, Xiu
717 Li, Jie Zhou, and Jiwen Lu. Ordinalclip: Learning rank
718 prompts for language-guided ordinal regression. In *Ad-*
719 *vances in Neural Information Processing Systems*, pages
720 35313–35325, 2022. 1, 2, 5, 6, 7
- 721 [38] Yuhong Li, Xiaofan Zhang, and Deming Chen. Csrnet: Di-
722 lated convolutional neural networks for understanding the
723 highly congested scenes. In *Proceedings of the IEEE Con-*
724 *ference on Computer Vision and Pattern Recognition*, pages
725 1091–1100, 2018. 2
- 726 [39] Feng Liang, Bichen Wu, Xiaoliang Dai, Kunpeng Li, Yinan
727 Zhao, Hang Zhang, Peizhao Zhang, Peter Vajda, and Diana
728 Marculescu. Open-vocabulary semantic segmentation with
mask-adapted clip. In *Proceedings of the IEEE/CVF con-*
ference on computer vision and pattern recognition, pages
7061–7070, 2023. 1
- [40] Hui Lin, Zhiheng Ma, Rongrong Ji, Yaowei Wang, and Xi-
aopeng Hong. Boosting crowd counting via multifaceted at-
tention. In *Proceedings of the IEEE/CVF conference on com-*
puter vision and pattern recognition, pages 19628–19637,
2022. 7
- [41] Wei Lin, Chenyang Zhao, and Antoni B Chan. Point-to-
region loss for semi-supervised point-based crowd counting.
In *Proceedings of the Computer Vision and Pattern Recogni-*
tion Conference, pages 29363–29373, 2025. 6, 7
- [42] Yanzhu Liu, Adams Wai Kin Kong, and Chi Keong Goh. A
constrained deep neural network for ordinal regression. In
Proceedings of the IEEE conference on computer vision and
pattern recognition, pages 831–839, 2018. 6, 7
- [43] Yanzhu Liu, Fan Wang, and Adams Wai Kin Kong. Proba-
bilistic deep ordinal regression based on gaussian processes.
In *Proceedings of the IEEE/CVF International Conference*
on Computer Vision, pages 5301–5309, 2019. 6
- [44] Zhiheng Ma, Xing Wei, Xiaopeng Hong, and Yihong Gong.
Bayesian loss for crowd count estimation with point super-
vision. In *Proceedings of the IEEE/CVF international con-*
ference on computer vision, pages 6142–6151, 2019. 5, 7
- [45] Siva Karthik Mustikovela, Varun Jampani, Shalini De Mello,
Sifei Liu, Umar Iqbal, Carsten Rother, and Jan Kautz. Self-
supervised viewpoint learning from image collections. In
Proceedings of the IEEE/CVF Conference on Computer Vi-
sion and Pattern Recognition, pages 3971–3981, 2020. 6
- [46] Thu Nguyen-Phuoc, Chuan Li, Lucas Theis, Christian
Richardt, and Yong-Liang Yang. Hologan: Unsupervised
learning of 3d representations from natural images. In *Pro-*
ceedings of the IEEE/CVF international conference on com-
puter vision, pages 7588–7597, 2019. 6
- [47] Zhenxing Niu, Mo Zhou, Le Wang, Xinbo Gao, and Gang
Hua. Ordinal regression with multiple output cnn for age
estimation. In *Proceedings of the IEEE Conference on Com-*
puter Vision and Pattern Recognition, pages 4920–4928,
2016. 1
- [48] Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy
Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez,
Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, et al.
Dinov2: Learning robust visual features without supervision.
arXiv preprint arXiv:2304.07193, 2023. 5, 6
- [49] Frank Palermo, James Hays, and Alexei A Efros. Dating his-
torical color images. In *European Conference on Computer*
Vision, pages 499–512. Springer, 2012. 1, 5
- [50] Hongyu Pan, Hu Han, Shiguang Shan, and Xilin Chen.
Mean-variance loss for deep age estimation from a face. In
Proceedings of the IEEE Conference on Computer Vision
and Pattern Recognition, pages 5285–5294, 2018. 2
- [51] Dileepa Pitawela, Gustavo Carneiro, and Hsiang-Ting Chen.
Cloc: Contrastive learning for ordinal classification with
multi-margin n-pair loss. In *Proceedings of the Computer*
Vision and Pattern Recognition Conference, pages 15538–
15548, 2025. 2, 5, 6, 7, 8
- [52] Monika Prokop-Piotrkowska, Kamila Marszałek-Dziuba,
Elżbieta Moszczyńska, Mieczysław Szalecki, and Elżbieta

787	Jurkiewicz. Traditional and new methods of bone age assessment-an overview. <i>Journal of clinical research in pediatric endocrinology</i> , 13(3):251, 2021. 1		
788			
789			
790	[53] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In <i>International Conference on Machine Learning</i> , pages 8748–8763. PMLR, 2021. 1, 2		
791			
792			
793			
794			
795			
796	[54] Yasiru Ranasinghe, Nithin Gopalakrishnan Nair, Wele Gedara Chaminda Bandara, and Vishal M Patel. Crowd-diff: Multi-hypothesis crowd density estimation using diffusion models. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pages 12809–12819, 2024. 1		
797			
798			
799			
800			
801			
802	[55] Karl Ricanek and Tamirat Tesafaye. Morph: A longitudinal image database of normal adult age-progression. In <i>7th international conference on automatic face and gesture recognition (FGRO6)</i> , pages 341–345. IEEE, 2006. 5		
803			
804			
805			
806	[56] Rasmus Rothe, Radu Timofte, and Luc Van Gool. Dex: Deep expectation of apparent age from a single image. In <i>Proceedings of the IEEE International Conference on Computer Vision Workshops</i> , pages 10–15, 2015. 1, 2		
807			
808			
809			
810	[57] Rasmus Rothe, Radu Timofte, and Luc Van Gool. Deep expectation of real and apparent age from a single image without facial landmarks. <i>International Journal of Computer Vision</i> , 126(2):144–157, 2018. 5		
811			
812			
813			
814	[58] Nataniel Ruiz, Eunji Chong, and James M Rehg. Fine-grained head pose estimation without keypoints. In <i>Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops</i> , pages 2074–2083, 2018. 1		
815			
816			
817			
818	[59] Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. Distilbert, a distilled version of bert: smaller, faster, cheaper and lighter. <i>arXiv preprint arXiv:1910.01108</i> , 2019. 5		
819			
820			
821			
822	[60] Rossano Schifanella, Miriam Redi, and Luca Maria Aiello. An image is worth more than a thousand favorites: Surfacing the hidden beauty of flickr pictures. In <i>Proceedings of the international AAAI conference on web and social media</i> , pages 397–406, 2015. 5		
823			
824			
825			
826			
827	[61] Wei Shen, Yilu Guo, Yan Wang, Kai Zhao, Bo Wang, and Alan L Yuille. Deep regression forests for age estimation. In <i>Proceedings of the IEEE conference on computer vision and pattern recognition</i> , pages 2304–2313, 2018. 5		
828			
829			
830			
831	[62] Nyeong-Ho Shin, Seon-Ho Lee, and Chang-Su Kim. Moving window regression: A novel approach to ordinal regression. In <i>Proceedings of the IEEE/CVF conference on computer vision and pattern recognition</i> , pages 18760–18769, 2022. 6		
832			
833			
834			
835	[63] Weibo Shu, Jia Wan, Kay Chen Tan, Sam Kwong, and Antoni B Chan. Crowd counting in the frequency domain. In <i>Proceedings of the IEEE/CVF conference on computer vision and pattern recognition</i> , pages 19618–19627, 2022. 7		
836			
837			
838			
839	[64] Vishwanath A Sindagi, Rajeev Yasarla, and Vishal M Patel. Jhu-crowd++: Large-scale crowd counting dataset and a benchmark method. <i>IEEE transactions on pattern analysis and machine intelligence</i> , 44(5):2594–2609, 2020. 5		
840			
841			
842			
843	[65] Qingyu Song, Changan Wang, Zhengkai Jiang, Yabiao Wang, Ying Tai, Chengjie Wang, Jilin Li, Feiyue Huang,	and Yang Wu. Rethinking counting and localization in crowds: A purely point-based framework. In <i>Proceedings of the IEEE/CVF international conference on computer vision</i> , pages 3365–3374, 2021. 6, 7	845 846 847 848
844			
	[66] Qingyu Song, Changan Wang, Yabiao Wang, Ying Tai, Chengjie Wang, Jilin Li, Jian Wu, and Jiayi Ma. To choose or to fuse? scale selection for crowd counting. In <i>Proceedings of the AAAI conference on artificial intelligence</i> , pages 2576–2583, 2021. 7		849 850 851 852 853
	[67] Satyajit Tourani, Siddharth Tourani, Arif Mahmood, and Muhammad Haris Khan. Unsupervised discovery of facial landmarks and head pose. In <i>Proceedings of the Computer Vision and Pattern Recognition Conference</i> , pages 21192–21202, 2025. 6		854 855 856 857 858
	[68] Roberto Valle, José M Buenaposada, and Luis Baumela. Multi-task head pose estimation in-the-wild. <i>IEEE Transactions on Pattern Analysis and Machine Intelligence</i> , 43(8):2874–2881, 2020. 6		859 860 861 862
	[69] Jia Wan and Antoni Chan. Adaptive density map generation for crowd counting. In <i>Proceedings of the IEEE/CVF international conference on computer vision</i> , pages 1130–1139, 2019. 1		863 864 865 866
	[70] Boyu Wang, Huidong Liu, Dimitris Samaras, and Minh Hoai Nguyen. Distribution matching for crowd counting. <i>Advances in neural information processing systems</i> , 33:1595–1607, 2020. 7		867 868 869 870
	[71] Jinhong Wang, Jintai Chen, Jian Liu, Dongqi Tang, Danny Z Chen, and Jian Wu. A survey on ordinal regression: Applications, advances and prospects. <i>arXiv preprint arXiv:2503.00952</i> , 2025. 2		871 872 873 874
	[72] Jinhong Wang, Shuo Tong, Jian Liu, Dongqi Tang, Weiqiang Wang, Wentong Li, Hongxia Xu, Danny Z Chen, Jintai Chen, and Jian Wu. Orderchain: Towards general instruct-tuning for stimulating the ordinal understanding ability of mllm. In <i>Proceedings of the IEEE/CVF International Conference on Computer Vision</i> , pages 3477–3487, 2025. 6		875 876 877 878 879 880
	[73] Rui Wang, Peipei Li, Huaibo Huang, Chunshui Cao, Ran He, and Zhaofeng He. Learning-to-rank meets language: Boosting language-driven ordering alignment for ordinal classification. In <i>Advances in Neural Information Processing Systems</i> , 2023. 1, 2, 5, 6, 7, 8		881 882 883 884 885
	[74] Xin Wang, Yudong Chen, and Wenwu Zhu. A survey on curriculum learning. <i>IEEE Transactions on Pattern Analysis and Machine Intelligence</i> , 44(9):4555–4576, 2022. 2, 3		886 887 888
	[75] Xinlong Wang, Xiaosong Zhang, Yue Cao, Wen Wang, Chunhua Shen, and Tiejun Huang. Seggpt: Towards segmenting everything in context. In <i>Proceedings of the IEEE/CVF International Conference on Computer Vision</i> , pages 1130–1140, 2023. 1		889 890 891 892 893
	[76] Yaoming Wang, Yangzhou Jiang, Jin Li, Bingbing Ni, Wenrui Dai, Chenglin Li, Hongkai Xiong, and Teng Li. Contrastive regression for domain adaptation on gaze estimation. In <i>Proceedings of the IEEE/CVF conference on computer vision and pattern recognition</i> , pages 19376–19385, 2022. 5		894 895 896 897 898
	[77] Daphna Weinshall, Gad Cohen, and Dan Amir. Curriculum learning by transfer learning: Theory and experiments with deep networks. In <i>International Conference on Machine Learning</i> , pages 5238–5246. PMLR, 2018. 3		899 900 901 902

903	[78] Tsun-Yi Yang, Yi-Ting Chen, Yen-Yu Lin, and Yung-Yu	
904	Chuang. Fsa-net: Learning fine-grained structure aggrega-	
905	tion for head pose estimation from a single image. In <i>Pro-</i>	
906	<i>ceedings of the IEEE/CVF Conference on Computer Vision</i>	
907	<i>and Pattern Recognition</i> , pages 1087–1096, 2019. 1, 6	
908	[79] Pengwei Yin, Jingjing Wang, Guanzhong Zeng, Di Xie, and	
909	Jiang Zhu. Lg-gaze: Learning geometry-aware continuous	
910	prompts for language-guided gaze estimation. In <i>European</i>	
911	<i>Conference on Computer Vision</i> , pages 1–17. Springer, 2024.	
912	1, 5	
913	[80] Pengwei Yin, Guanzhong Zeng, Jingjing Wang, and Di	
914	Xie. Clip-gaze: Towards general gaze estimation via visual-	
915	linguistic model. In <i>Proceedings of the AAAI Conference on</i>	
916	<i>Artificial Intelligence</i> , pages 6729–6737, 2024. 5	
917	[81] Boyang Yu, Aakash Kaku, Kangning Liu, Avinash Parnandi,	
918	Emily Fokas, Anita Venkatesan, Natasha Pandit, Rajesh	
919	Ranganath, Heidi Schambra, and Carlos Fernandez-Granda.	
920	Quantifying impairment and disease severity using ai models	
921	trained on healthy subjects. <i>NPJ Digital Medicine</i> , 7(1):180,	
922	2024. 1	
923	[82] Jure Zbontar, Li Jing, Ishan Misra, Yann LeCun, and	
924	Stéphane Deny. Barlow twins: Self-supervised learning via	
925	redundancy reduction. In <i>International Conference on Ma-</i>	
926	<i>chine Learning</i> , pages 12310–12320. PMLR, 2021. 3	
927	[83] Kaiwen Zha, Peng Cao, Jeany Son, Yuzhe Yang, and Dina	
928	Katabi. Rank-n-contrast: learning continuous representa-	
929	tions for regression. <i>Advances in Neural Information Pro-</i>	
930	<i>cessing Systems</i> , 36:17882–17903, 2023. 2, 6	
931	[84] Cheng Zhang, Hai Liu, Yongjian Deng, Bochen Xie, and	
932	Youfu Li. Tokenhpe: Learning orientation tokens for effi-	
933	cient head pose estimation via transformers. In <i>Proceedings</i>	
934	<i>of the IEEE/CVF conference on computer vision and pattern</i>	
935	<i>recognition</i> , pages 8897–8906, 2023. 1	
936	[85] Xucong Zhang, Yusuke Sugano, Mario Fritz, and Andreas	
937	Bulling. Appearance-based gaze estimation in the wild. In	
938	<i>Proceedings of the IEEE Conference on Computer Vision</i>	
939	<i>and Pattern Recognition</i> , pages 4511–4520, 2015. 1, 2	
940	[86] Xucong Zhang, Yusuke Sugano, Mario Fritz, and Andreas	
941	Bulling. Mpiigaze: Real-world dataset and deep appearance-	
942	based gaze estimation. <i>IEEE transactions on pattern analy-</i>	
943	<i>sis and machine intelligence</i> , 41(1):162–175, 2017. 5	
944	[87] Xucong Zhang, Seonwook Park, Thabo Beeler, Derek	
945	Bradley, Siyu Tang, and Otmar Hilliges. Eth-xgaze: A large	
946	scale dataset for gaze estimation under extreme head pose	
947	and gaze variation. In <i>European conference on computer vi-</i>	
948	<i>sion</i> , pages 365–381. Springer, 2020. 5	
949	[88] Yingying Zhang, Desen Zhou, Siqin Chen, Shenghua Gao,	
950	and Yi Ma. Single-image crowd counting via multi-column	
951	convolutional neural network. In <i>Proceedings of the IEEE</i>	
952	<i>Conference on Computer Vision and Pattern Recognition</i> ,	
953	pages 589–597, 2016. 5	
954	[89] Wenqi Zhong, Chen Xia, Dingwen Zhang, and Junwei Han.	
955	Uncertainty modeling for gaze estimation. <i>IEEE Transac-</i>	
956	<i>tions on Image Processing</i> , 33:2851–2866, 2024. 1	
957	[90] Kaiyang Zhou, Jingkang Yang, Chen Change Loy, and Ziwei	
958	Liu. Learning to prompt for vision-language models. <i>In-</i>	
959	<i>ternational Journal of Computer Vision</i> , 130(9):2337–2348,	
960	2022. 5, 6, 7	
	[91] Yi Zhou, Connelly Barnes, Jingwan Lu, Jimei Yang, and Hao	961
	Li. On the continuity of rotation representations in neural	962
	networks. In <i>Proceedings of the IEEE/CVF conference on</i>	963
	<i>computer vision and pattern recognition</i> , pages 5745–5753,	964
	2019. 2	965
	[92] Xiangyu Zhu, Zhen Lei, Junjie Yan, Dong Yi, and Stan Z	966
	Li. High-fidelity pose and expression normalization for face	967
	recognition in the wild. In <i>Proceedings of the IEEE con-</i>	968
	<i>ference on computer vision and pattern recognition</i> , pages	969
	787–796, 2015. 5	970
	[93] Xiangyu Zhu, Zhen Lei, Xiaoming Liu, Hailin Shi, and	971
	Stan Z Li. Face alignment across large poses: A 3d solu-	972
	tion. In <i>Proceedings of the IEEE conference on computer</i>	973
	<i>vision and pattern recognition</i> , pages 146–155, 2016. 5	974