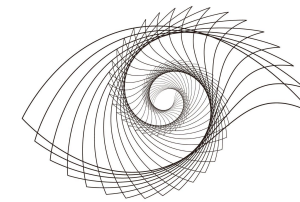# CVPR 2025 Report

Hirokatsu Kataoka, Yoshihiro Fukuhara,

Ryousuke Yamada, Daichi Otsuka, Rintaro Yanagi, Oishi Deb,
Kazuya Nishimura, Moeri Okuda, Yuto Matsuo, Ren Ohkubo, Yue Qiu, Noritake Kodama,
Gido Kato, Kenzo Yamabuki, Joe Hasei, Ryuichi Nakahara, Yukinori Yamamoto, Sho
Okazaki, Kohsuke Ide, Yuiga Wada,
Daichi Yashima, Shinichi Mae, Hinako Mitsuoka, Maika Takada,
Jianyuan Wang

LIMIT.Lab / cvpaper.challenge / Visual Geometry Group (VGG)

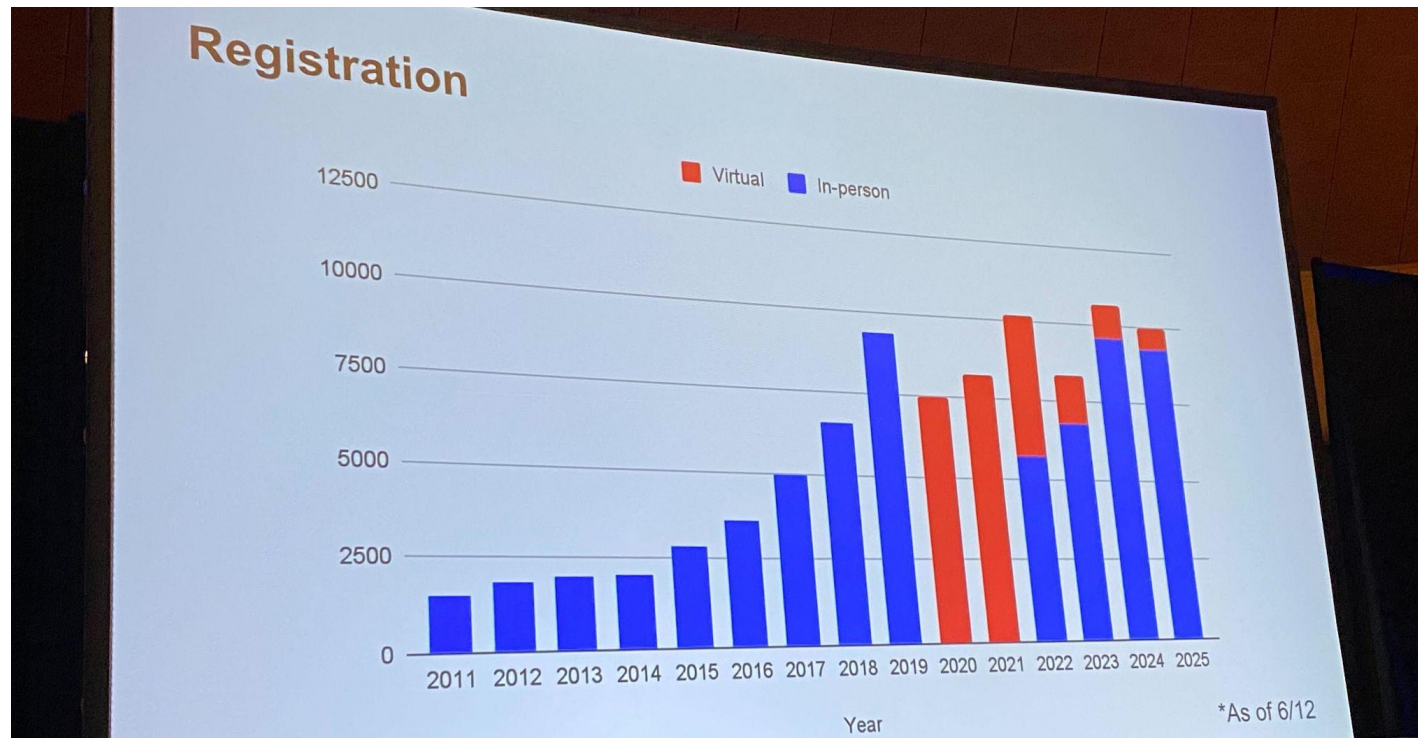# Meta Insights into Trends and Tendencies
# in CVPR 2025

- What kind of research is currently trending?

- What are overseas researchers working on?

- We have compiled the "trends" and "insights."

## From opening slide

- ❏ The attendees slightly decreased this year
  - ❏ By comparing to the CVPR 2024 attendees
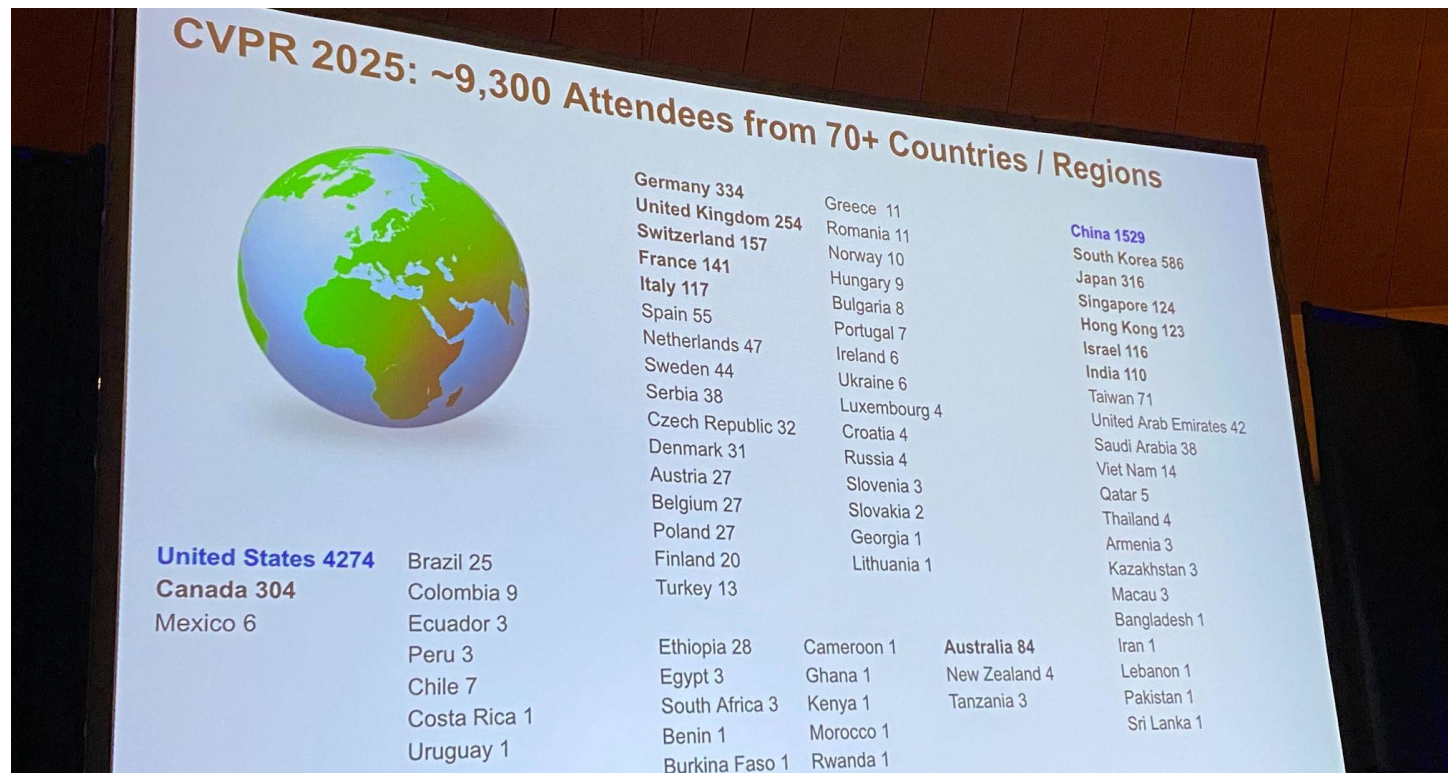- ❏ Virtual attendance is being allowed

## From opening slide

❏ Listed at each continent
- ❏ US: 4,274, China: 1,529, South Korea: 586, Germany: 334, Japan: 316, Canada: 304, United Kingdom: 254
- ❏ Totally 70+ countries/regions

## From opening slide

- ❏ The #workshops slightly decreased from the previous year
  - ❏ 2024: 123 vs. 2025: 118 workshops
- ❏ The 118 ws have been divided into 26 tracks

## From opening slide

- ❏ You can see the leading companies in CVPR from the list!
- ❏ Thanks to the companies, the CVPR community has been accelerated

## From opening slide

- ❏ #Submissions − 2024: 11,532 vs. 2025: 13,008 (12.8% increased)
- ❏ Acceptance rate − 2024: 23.55% vs. 2025: 22.08% (1.47 pt decreased)

## From opening slide

❏ Rapidly growing CVPR community!
  ❏ in terms of #authors, #reviewers, and #area chairs

## From opening slide

❑ Newly induced CVPR paper from reviewer contribution

    ❑ A paper included outstanding reviewer is highlighted

    ❑ Review contributions may influence accepted papers!

## From opening slide

❏ Listed the recent trends in CVPR

    ❏ What are the next trends?

## From opening slide

❏ The program committee is struggling to maintain the review quality
  ❏ Some decisions for reviewing process
  ❏ Qualified reviewers / reviewer education should be required

## From opening slide

❏ 1st author / multi-paper submitter must review others

  ❏ It's now mandatory

  ❏ Only voluntary review is difficult

## From opening slide

- ❏ According to the area chairs,
    - ❏ Relatively better reviews have been gathered in CVPR 2025
    - ❏ A Ph.D. student can be a good reviewer!

## From award ceremony

## Summary of best paper candidates 1: Molmo and PixMo: Open Weights and Open Data for State-of- the-Art Vision-Language Models

❏ Molmo is a fully open VLM family built without using proprietary data or models, addressing the lack of foundational knowledge in open VLM development.

❏ Trained on the newly collected PixMo datasets, Molmo's 72B model outperforms major proprietary models like Gemini 1.5 Pro and ranks just behind GPT-4o.

## Summary of best paper candidates 2: Foundation Stereo: Zero-shot Stereo Matching

❏ Foundation Stereo is a foundation model for stereo depth estimation, achieving strong zero-shot generalization without target-domain fine-tuning

❏ It is trained on a large-scale synthetic dataset and enhanced by self-curation, monocular prior adaptation, and long-range context reasoning for cross-domain robustness

## Summary of best paper candidates 3: VGGT: Visual Geometry Grounded Transformer

❏ VGGT is a feed-forward network that predicts full 3D scene attributes from one or more views, including depth, camera parameters, and 3D point tracks.

❏ It achieves state-of-the-art performance across multiple 3D tasks with high speed and efficiency, and boosts downstream tasks when used as a feature backbone.

## Summary of best paper candidates 4: Descriptor-In-Pixel: Point-Feature Tracking For Pixel Processor Arrays

❑ This paper introduces a fully in-pixel method for point-feature detection and tracking on PPA vision sensors, eliminating the need to output raw image data.

❑ Using a Descriptor-In-Pixel paradigm, the system runs at over 3000 FPS on the SCAMP-7 prototype, achieving over 1000× data reduction while maintaining reliable tracking under fast motion.

## Summary of best paper candidates 5: The PanAf-FGBG Dataset: Understanding the Impact of Backgrounds in Wildlife Behaviour Recognition

- ❏ PanAf-FGBG is a novel dataset of wild chimpanzee behaviours paired with background videos from the same camera locations, enabling direct evaluation of background impact on behaviour recognition.

- ❏ It supports in- and out-of-distribution testing and introduces a normalization technique that significantly improves model generalization.

Summary of best paper candidates 6: MegaSaM: Accurate, Fast and Robust Structure and

Motion from Casual Dynamic Videos

- ❏ This work presents a deep visual SLAM system that accurately estimates camera poses and depth maps from monocular videos of dynamic scenes with minimal parallax.
- ❏ With tailored training and inference schemes, the system outperforms existing methods in both accuracy and speed on real and synthetic datasets.



LIMIT.LAB

https://limitlab.xyz/

## Summary of best paper candidates 7: TacoDepth: Towards Efficient Radar-Camera Depth Estimation with One-stage Fusion

❏ TacoDepth is a one-stage Radar-Camera depth estimation model that improves efficiency and robustness without relying on intermediate depth results.

❏ It achieves 12.8% better accuracy and 91.8% faster processing than previous methods, offering a strong balance between speed and performance for real-time applications.

## Summary of best paper candidates 8: Navigation World Models

❏ NWM is a controllable video generation model that predicts future views for navigation using a Conditional Diffusion Transformer trained on egocentric videos.

❏ It enables dynamic, constraint-aware planning and can simulate trajectories even in unfamiliar environments from a single image.

## Summary of best paper candidates 9: Convex Relaxation for Robust Vanishing Point Estimation in Manhattan World

❏ GlobustVP introduces a convex relaxation approach for vanishing point detection in Manhattan worlds, combining soft line–VP associations with semidefinite programming.

❏ It achieves global optimality with high efficiency and robustness, outperforming prior methods on both synthetic and real–world datasets.



Ground Truth      GlobustVP (Ours)

3 VPs
38 lines

100%, 100%
0.23 pix., 50 ms

LIMIT.LAB
https://limitlab.xyz/

22

## Summary of best paper candidates 10: UniAP: Unifying Inter- and Intra-Layer Automatic Parallelism by Mixed Integer Quadratic Programming

❑ UniAP is the first automatic parallelism method that jointly optimizes inter- and intra-layer strategies using mixed integer quadratic programming.

❑ It achieves up to 3.80× higher throughput and 107× faster optimization than previous state-of-the-art methods on Transformer-based models.

## Summary of best paper candidates 11: Zero-Shot Monocular Scene Flow Estimation in the Wild

❏ This work introduces a generalizable scene flow model that jointly estimates geometry and motion, addressing key limitations in current approaches.

❏ It leverages a diverse synthetic dataset of 1M samples and achieves strong zero-shot performance on real-world videos and robotic scenes.

## Summary of best paper candidates 12: 3D Student Splatting and Scooping

❏ This paper introduces Student Splatting and Scooping (SSS), a new mixture model using Student's t distributions with both positive and negative densities for novel view synthesis.

❏ SSS improves quality and parameter efficiency over standard 3D Gaussian Splatting, reducing components by up to 82% while maintaining or surpassing rendering performance.



(a) Bonsai from Mip-NeRF 360 [2]
(b) Gardern from Mip-NeRF 360 [2]
(c) Kitchen from Mip-NeRF 360 [2]
(d) Truck from Tanks & Temples [16]
(e) Playroom from Deep Blending [10]



Comparison of Student's t-Distributions with Various Degrees of Freedom

## Summary of best paper candidates 13: DIFIX3D+: Improving 3D Reconstructions with Single-Step Diffusion Models

❏ DIFIX3D+ is a novel pipeline that enhances 3D reconstruction and novel-view synthesis using a single-step diffusion model called DIFIX to remove artifacts.

❏ Compatible with both NeRF and 3DGS, it improves underconstrained regions and achieves 2× better FID scores while preserving 3D consistency.

## Summary of best paper candidates 14: Generative Multimodal Pretraining with Discrete Diffusion Time-step Tokens

❏ This work introduces a unified framework that combines LLMs and diffusion models using recursive visual tokens derived from diffusion timesteps.

❏ The approach enables effective multimodal comprehension and generation, outperforming existing MLLMs across both tasks.

## From award ceremony



Reference: CVPR X website

## From award ceremony



*Reference: CVPR X website*

## From award ceremony



*Reference: CVPR X website*

## From award ceremony

## From award ceremony

From award ceremony

## From award ceremony

## From award ceremony

## From award ceremony

## From award ceremony

From award ceremony

LIMIT.LAB
https://limitlab.xyz/

From award ceremony

## From award ceremony

From award ceremony

## From award ceremony

## From award ceremony

## Unofficial best paper prediction in LIMIT.Lab

❏ This initiative was organized by a group of volunteers to coincide with the release of the Best Paper Award Candidates for CVPR 2025.

❏ Independent from the official review process, it aims to predict which paper will receive the award, based on our own perspectives.

❏ Beyond simple prediction, the goal is to deepen our understanding of the field by carefully reading the nominated papers and evaluating their novelty, impact, and future potential.

❏ Through discussion and exchanging opinions, the initiative also helps broaden our perspectives and refine our criteria for assessing research.

LIMIT.LAB

https://limitlab.xyz/

## Unofficial best paper prediction in LIMIT.Lab

- ❏ Our Award Selection Process
    - ❏ We initiated a call on Slack to form the (unofficial) CVPR 2025 Award Committee and created a dedicated Slack channel for coordination.
    - ❏ For every paper, at least one committee member read it thoroughly and created a summary. All members reviewed these summaries to ensure they had at least a basic understanding of every paper before forming their judgments.
    - ❏ Each member shared their individual award list, which was then discussed and consolidated through committee deliberation.
    - ❏ After forming the award list, the selected papers were re-reviewed before making the final decision.
- ❏ Selection Criteria
    - ❏ The evaluation primarily follows CVPR regulations, focusing on the level of contribution to the field and the potential to significantly impact its future.
    - ❏ Of course, factors such as author prominence or existing citation count are excluded from consideration.

LIMIT.LAB 45

## Unofficial best paper prediction in LIMIT.Lab

- ❏ Our selected papers for each award
    - ❏ Total: 9 papers (referring the CVPR 2024 list; 10 awarded papers in CVPR 2024)
    - ❏ Best Student Paper Honorable Mention(2)
    - ❏ Best Student Paper Award(2)
    - ❏ Best Paper Honorable Mention(3)
    - ❏ Best Paper Award(2)
    - ❏ Note that the award candidate list ( https://cvpr.thecvf.com/virtual/2025/events/AwardCandidates2025 ) and the awardees are slightly different

## Unofficial best paper prediction in LIMIT.Lab

- ❏ Our selected papers for each award
  - ❏ Total: 9 papers from the CVPR 2024 list (10 awarded papers in CVPR 2024)
  - ❏ Best Student Paper Honorable Mention（2）
    - ❏ DIFIX3D+: Improving 3D Reconstructions with Single-Step Diffusion Models
    - ❏ 3D Student Splatting and Scooping
  - ❏ Best Student Paper Award（2）
    - ❏ Zero-Shot Monocular Scene Flow Estimation in the Wild
    - ❏ Convex Relaxation for Robust Vanishing Point Estimation in Manhattan World
  - ❏ Best Paper Honorable Mention（3）
    - ❏ Navigation World Models
    - ❏ MegaSaM: Accurate, Fast and Robust Structure and Motion from Casual Dynamic Videos
    - ❏ Molmo and PixMo: Open Weights and Open Data for State-of-the-Art Vision-Language Models
  - ❏ Best Paper Award（2）
    - ❏ Generative Multimodal Pretraining with Discrete Diffusion Time-step Tokens
    - ❏ VGGT: Visual Geometry Grounded Transformer

## Unofficial best paper prediction in LIMIT.Lab

- ❏ Our selected papers for each award <mark>and reality</mark>
  - ❏ Total: 9 → <mark>7</mark> papers from the CVPR 2024 <u>list</u> (10 awarded papers in CVPR 2024)
  - ❏ Best Student Paper Honorable Mention（2→ <mark>1</mark>）
    - ❏ ❌DIFIX3D+: Improving 3D Reconstructions with Single-Step Diffusion Models
    - ❏ ✅ 3D Student Splatting and Scooping → But, it's <mark>Best Paper Honorable Mention</mark>
  - ❏ Best Student Paper Award（2→ <mark>1</mark>）
    - ❏ ❌Zero-Shot Monocular Scene Flow Estimation in the Wild
    - ❏ ❌Convex Relaxation for Robust Vanishing Point Estimation in Manhattan World
  - ❏ Best Paper Honorable Mention（3→ <mark>4</mark>）
    - ❏ ✅ Navigation World Models
    - ❏ ✅ MegaSaM: Accurate, Fast and Robust Structure and Motion from Casual Dynamic Videos
    - ❏ ✅ Molmo and PixMo: Open Weights and Open Data for State-of-the-Art Vision-Language Models
  - ❏ Best Paper Award（2→ <mark>1</mark>）
    - ❏ ✅ Generative Multimodal Pretraining with Discrete Diffusion Time-step Tokens → But, it's <mark>Best Student Paper Honorable Mention</mark>
    - ❏ ✅ VGGT: Visual Geometry Grounded Transformer

## Meta-level insights and BP selected discussion points

❑ Even before reading the papers, the titles alone were fairly accurate, 6 out of the 7, picked matched well. This shows that it's important to craft titles that clearly convey the contribution and impact to some extent.

❑ You may not win an award even if you aim for it, but you'll never reach it if you don't aim at all!

❑ VGGT is the result of years of accumulated effort, including extensive method exploration and parameter tuning, as well as close collaboration within the lab, combining components like VGGSfM and CoTracker.

## Meta-level insights and BP selected discussion points

- ❏ 3D reconstruction was once again a major highlight this year.
- ❏ 3D Reconstruction → Faster and higher-quality results!
- ❏ VLMs (Vision-Language Models) → Open-source momentum and improved controllability over images!
- ❏ Robotics → Groundbreaking applications of world models!

## Meta-level insights and BP selected discussion points

❏ Reading the CVPR Best Paper candidate studies felt like diving into a collection of science fiction stories, each one filled with a sense of wonder and possibility.

❏ These works weren't just about incremental technical improvements; they felt like presentations of futuristic tools. They made me think, "If something this amazing is possible, I'd want to build something like this." Each paper sparked real inspiration.

❏ What made them even more compelling was how thoroughly the authors validated their ideas through meticulous experimentation.

❏ Just like good science fiction is grounded in a sense of realism, these carefully designed experiments gave the research a tangible, exciting credibility.

## Backstory of the Best Paper Award – VGGT (Thank you, Jianyuan Wang, for sharing the story!)

- ❏ **VGGSfM in the beginning**
  - ❏ Was too complex and difficult to get started with
  - ❏ Follows the classical SfM pipeline: images – correspondences – camera poses & points – bundle adjustment – refined camera poses & points
- ❏ **VGGT authors also noticed DUSt3R**
  - ❏ Reveals the power of large-scale data-driven paradigm
  - ❏ Went for VGGT – Transformer feedforward only, everyone can use it
- ❏ **VGG (Lab) has proposed CoTracker for an accurate point tracker**



Figure 2. **Overview of VGGSfM.** Our method extracts 2D tracks from input images, reconstructs cameras using image and track features, initializes a point cloud based on these tracks and camera parameters, and applies a bundle adjustment layer for reconstruction refinement. The whole framework is fully differentiable and designed for end-to-end training.

J. Wang et al. "VGGSfM: Visual Geometry Grounded Deep Structure From Motion," CVPR 2024. [Link]



Fig. 3: **CoTracker architecture.** We compute convolutional features $\phi(I_t)$ for every frame and process them with sliding windows. To initialize track features $Q$, we bilinearly sample from $\phi(I_t)$ with starting point locations $P$. Locations $P$ also serve to initialize estimated tracks $\hat{P}$. See Fig. 4 for a visualization of one sliding window.

N. Karaev et al. "CoTracker: It is Better to Track Together," ECCV 2024. [Link]

## Backstory of the Best Paper Award – VGGT <span>(Thank you, Jianyuan Wang, for sharing the story!)</span>

❏ Flight from 🇬🇧 to 🇺🇸

  ❏ Has been waiting for the passport

  ❏ Got his passport just before the departure

  ❏ Arrived Nashville at midnight before the award ceremony

❏ Award ceremony

  ❏ No contact from the Program Chairs about the award, CVPR this year felt more like the Oscars

  ❏ Best Paper Award was saved for last (this means it was 'all or nothing' after the honorable mention)

  ❏ Comments from Jianyuan:

  I was extremely nervous. I still don't know how to fully describe the feeling.

  Even now, it all feels kind of unreal.

## Backstory of the Best Paper Award – VGGT

- ❏ Perfect switching from VGGNet to VGGTransformer!
  - ❏ VGGNet (ICLR 2015)
    - ❏ Convolutional Neural Networks (CNN)
    - ❏ A backbone for image recognition
  - ❏ VGGTransformer / VGGT (CVPR 2025) – ten years after the VGGNet
    - ❏ Transformer
    - ❏ 3D reconstruction in a single Transformer
    - ❏ A backbone for 3D vision
  - ❏ Rough history of Visual Geometry Group (VGG)
    - ❏ Before 2010s: 3D geometry is the main focus, but there are recognition researches
    - ❏ After 2010s:   Deep learning is the main focus, but there are 3D vision researches
    - ❏ VGGT: The strong cross point between 3D geometry & deep learning
    - ❏ After 2025?: Everyone is able to participate in 3D x deep learning?

## Trends of Medical x CV research

a. Vision and language for medical domain

    i. Medical VQA with noisy labels and diffusion [Guo+, CVPR 2025]

    ii. VLMs alignment with CoOP [Koleilat+, CVPR 2025]

    iii. VQA with visual reference [Chen+, CVPR 2025]

    iv. VLM with soft label [Ko+, CVPR 2025]

    v. Pre-training of VLM [Ziyang+, CVPR2025]

b. Semi-supervised learning for medical

    i. The task has been focused on for a long time, but it is still prevalent this year.

    ii. Is this evidence that the lack of labeled data has not been resolved, even with the foundation model in medical?

    iii. Deal with annotation ambiguity [Kumari+, CVPR 2025]

    iv. Depth guided segmentation [Li+, CVPR 2025]

    v. Find overconfidence prediction of foundation model [Ma+, CVPR 2025]

    vi. Unsupervised prompting for SAM by DPO-inspired loss [Konwer+, CVPR 2025]

    vii. Uncertainty-aware consistency and contrastive [Assefa+, CVPR 2025]

*LIMIT.LAB*

https://limitlab.xyz/

## Trends of Pathology x CV research

a. Vision and language for pathology

    i. Patch-level VLLM -> multi-resolution (including patch, slide-level) VLLM

[Chen+, CVPR 2025]  [Sun+, CVPR 2025]  [Albastaki+, CVPR 2025]

b. mamba-based model for whole slide image

    i. To effectively aggregate patch-level information for mamba

[Zhang+, CVPR 2025]  [Zheng+, CVPR 2025]

Figure 1. **Left: Conventional MIL Bagging** of patches adopts *no spatial context*. **Middle: 1D Mamba-based methods** flatten a WSI into a 1D sequence and lose the 2D structure. The adjacent blue and orange patches are far away in the sequence. We call this "*spatial discrepancy*". **Right: 2DMamba** processes a WSI in a 2D manner, preserving 2D structure. The blue and orange patches maintains adjacent in the sequence. We call it "*spatial continuity*".

LIMIT.LAB

## Direct Preference optimization (DPO) for Vision and language

### Boosting for in-context leaning with SymDPO
[Jia+, cvpr 2025]



Figure 2. Comparison of General DPO and SymDPO Formats: General DPO relies solely on standard text for Questions, Answers, Chosen, and Rejected Answers, focusing on text-based training. In contrast, SymDPO replaces textual Answers with symbolized text to boost multimodal understanding, requiring models to interpret both visual and symbolized cues. This approach strengthens the model's ability to reason and decide in complex multimodal contexts.

### Debiasing VLM with Noise-aware preference optimization
[Zhang+, cvpr 2025]



Figure 2. **Method details.** First, biased responses are constructed by using masking to guide the model toward over-relying on prompts and generating responses based on the base model. Next, NaPO is applied for noise-robust preference optimization to counteract noise in automatically constructed data, dynamically assessing data noise levels to calculate NaPO's noise robustness coefficient $q$ (see Equation (12)). Here we assumed that the original data is of high quality, so DPO is used to train on it directly. Additional experiments were conducted with NaPO on the original data, and the results can be found in Appendix A.

### Improving VLM with task preference optimization by introduce task specific token
[Yan+, cvpr 2025]



Figure 2. **Comparison of Learning Method.** A solid line indicates data flow, and a dotted line represents feedback. ❄ and 🔥 denote modules that are frozen and unfrozen.



Figure 3. **Overall Pipeline of TPO.** The architecture of Task Preference Optimization (TPO) consists of four main components: (1) a vision encoder, (2) a connector, (3) a large language model, and (4) a series of visual task heads. Differently colored flame symbols indicate which components are unfrozen at various stages of the training process.

## Direct Preference optimization (DPO) for diffusion model

### Step-by-step preference optimization
### use PO for each step [Liang+, cvpr 2025]



Figure 3. Comparing frameworks of SPO, Diffusion-DPO, and D3PO approaches. SPO does not adopt direct preference propagation as other DPO methods do. In SPO, a pool of samples are generated at each step, from which a proper win/lose pair is selected and used to fine-tune the diffusion model. Then, a single sample is randomly selected to initialize the next iteration.

### DPO with curriculum learning
### [NA+, cvpr 2025] [Croitoru+, cvpr 2025]



Figure 2. **Three-stage training of HG-DPO.** It progressively enhances the model's human image generation capabilities.

### Calibrated Preference Optimization
### using multiple reward model [Lee+, cvpr 2025]



Figure 2. **Overview.** (a) We generate $N$ images using pretrained T2I diffusion model using the prompt dataset, and infer the scores from reward models. (b) Then, we calibrate the rewards by making pairwise comparison between images. For each image, we compute the win-rates between other $N-1$ images using Eq. (2), and average them to obtain calibrated reward $R_{ca}$ (see Sec. 4.2). (c) We select pair by choosing the best-of-$N$ and worst-of-$N$ when using single reward. For multi-reward, we use non-dominated sorting algorithm to select upper Pareto set as positives, and lower Pareto set as negatives. The accepted and rejected pairs are also listed using proposed rejection sampling method. (d) Lastly, during training, we select a pair from (c), and compute CaPO loss (i.e., Eq. (8)), which perform regression task to match the difference in calibrated rewards (i.e., $\Delta R_{ca}$ by the difference of implicit reward model (i.e., $\Delta R_\theta$).

### Inversion preference optimization
### with reparametrization DDIM [Lu+, cvpr 2025]



Figure 2. Illustration of Inversion for Preference Optimization.

## Direct Preference optimization (DPO) for downstream task
## May DPO-based fine-tuning become the next trend?

Enhancing SAM with Efficient Prompting and Preference Optimization
for Semi-supervised Medical Image Segmentation
[Konwer+, cvpr 2025]

- Efficient unsupervised prompting strategy
  that enhances segmentation performance
- Step1: Add preferences ratings for candidates
  that is estimated by vision and language model.
- Step2: Train model with DPO

We should tackle label-efficient dataset generation for DPO :
- Labor efficiency
  - Utilize foundation model?
- Ensure quality of data
  - Rethinking label-efficient techniques?
- Bias of data
  - Integrates multiple models?



Figure 2. **Illustration of the proposed framework for semi-supervised segmentation:** Unsupervised geometric and text prompts, obtained from pretrained BiomedCLIP, MedVInT, and GPT-4 models, are fed into the prompt encoder for finetuning the framework on a small fraction of annotated data. In the next stage, we simulate a virtual annotation process that assigns ratings to the generated segmentation candidates, which are used to fine-tune the decoder. This stage handles unannotated data, as the model does not rely on ground truth for direct supervision but only for rating while simulating a human annotator's feedback.

## Mitigating Hallucinations in Large Vision-Language Models via DPO: On-Policy Data Hold the Key

- ○ DPO cannot learn strictly off-policy preferred answers (reverse-KL → ∞)
- ○ Make expert-corrected responses on-policy before DPO
- ○ Overall Pipeline:
  - 1. Collect & Rate
    - Base LVLM → response candidates
    - GPT-4V tags sentence-level hallucination severity / error type and minimally revises text (4.8k pairs)
  - 2. On-Policy Alignment (OPA)
    - LoRA-SFT on original + revised answers
    - → brings corrections into model's support.
  - 3. OPA-DPO Fine-tune with three preference pairs
    - Language Correction (hallucination-weighted)
    - Image Focus (clean vs 30 %-masked image)
    - Anchored Preference (keep preferred probs from drifting)

## Neural Hierarchical Decomposition for Single Image Plant Modeling

❏ Abstract:

  ❏ Constructing high-quality 3D models of biological plants is challenging. They tackled on the 3D decomposition of plants from a single image.

  ❏ segmentation → structure inference → box decomposition → VAE

❏ Key point:

  ❏ This approach is adaptable for both of the indoor and outdoor trees



Zhihao et al., "Neural Hierarchical Decomposition for Single Image Plant Modeling", CVPR 2025, [Link]

LIMIT.LAB

https://limitlab.xyz/

## Synthetic data has high potential

❏ From the perspective of best paper candidates,

  ❏ 1. Synthetic data – Data generation using world simulator, generative models, and other resources

  ❏ 2. Language-free vision foundation model – Depth, optical flow, matting, segmentation

  ❏ 3. Zero-shot / in-the-wild recognition – Without an additional fine-tuning, the synthetic pre-training is working on real-world data



B. Wen et al. "FoundationStereo: Zero-Shot Stereo Matching," CVPR 2025. [Link]



Y. Liang et al. "Zero-Shot Monocular Scene Flow Estimation in the Wild," CVPR 2025. [Link]

LIMIT.LAB

https://limitlab.xyz/

## Synthetic data has high potential

❏ From the perspective of CVPR 2025 papers,

❏ Infinigen / Infinigen Indoors have been applied the following papers

  ❏ DI-PCG: An inverse PCG (Procedural Content Generation) framework that realizes image-to-3D generation via parameter estimation in diffusion models

  ❏ BlenderGym: A comprehensive benchmark for real-world graphics editing tasks by having VLMs to reconstruct 3D scenes from a start to a goal state.



W. Zhao et al. "DI-PCG: Diffusion-based Efficient Inverse Procedural Content Generation for High-quality 3D Asset Creation," CVPR 2025. [Link]

Y. Gu et al. "BlenderGym: Benchmarking Foundational Model Systems for Graphics Editing," CVPR 2025. [Link]

## Synthetic data has high potential

❏ From the perspective of CVPR 2025 SynData4CV Workshop,

❏ Synthetic / generated data can be used training for vision / VL models

  ❏ Over 60+ posters in the SynData4CV Workshop

  ❏ Many tasks such as classification, detection, segmentation, 3D recognition, anomaly detection, video, robot manipulation

The workshop aims to explore the use of synthetic data in training and evaluating computer vision models, as well as in other related domains. During the last decade, advancements in computer vision were catalyzed by the release of painstakingly curated human-labeled datasets. Recently, people have increasingly resorted to synthetic data as an alternative to laborintensive human-labeled datasets for its scalability, customizability, and costeffectiveness. Synthetic data offers the potential to generate large volumes of diverse and high-quality vision data, tailored to specific scenarios and edge cases that are hard to capture in real-world data. However, challenges such as the domain gap between synthetic and real-world data, potential biases in synthetic generation, and ensuring the generalizability of models trained on synthetic data remain. We hope the workshop can provide a forum to discuss and encourage further exploration in these areas.

Workshop overview



Best Long Paper: Synthetic data for robot policy learning

Z. Xue et al. "DemoGen: Synthetic Demonstration Generation for Data-Efficient Visuomotor Policy Learning," CVPR 2025 SynData4CV Workshop. [Link]

LIMIT.LAB

64

https://limitlab.xyz/

## Molmo & PixMo: State-of-the-art open-sourced VLM

❏ Molmo & PixMo have been proven:

  ❏ Data quality rather than simple model scale (once again!)

  ❏ Molmo keeps the standard architecture and simple VL connections

  ❏ For the purpose of open strategy, PixMo consists of public & synthetic data

## ViT (Vision Transformer) has been improved several points!

❏ Token Cropr ViT

    ❏ Token pruning that uses auxiliary prediction heads

❏ Your ViT is Secretly an Image Segmentation Model

    ❏ With learnable queries + mask logits, simple ViT can be a segmentor

## ViT (Vision Transformer) has been improved several points!

- ❏ BOE-ViT (Boosting Orientation Estimation ViT)
  - ❏ ViT for 3D subtomogram alignment through shift/rotation estimation
- ❏ BHViT (Binarized Hybrid ViT)
  - ❏ Full binarized ViT model / Xnor and popcount operations



…and some more methods regarding ViT architecture

## Synthetic data for computer vision

- ❏ **CLIPasso: Semantically-Aware Object Sketching**
- ❏ Generation Sketch images tasks is important.
  - ❏ Many artistic and scientific ideas begin as sketches.
- ❏ Challenges in generating sketch images by diffusion models.
  - ❏ Key point: collecting a large-scale sketch dataset is challenging because existing sketches are limited in number and often subject to copyright restrictions.
- ❏ Approach: Generating sketch images from large real-image datasets with ControlNet.
  - ❏ The synthetic sketches can then be used to train new diffusion models, enabling large-scale, copyright-free sketch generation.
- ❏ https://clipasso.github.io/clipasso/

## Graph structures for video understanding using VLMs

- ❏ [1] Spatio-temporal scene graph, enabling VLM to self-generate detailed multi-step reasoning data and improve their compositional reasoning abilities. (Qui et al.[Link],)

- ❏ [2]Topological semantic graph, empowering VLM to deeply understand the 3D context of complex long-form videos and perform human-like reasoning. (Huang et al.[Link])

- ❏ [3]Graph structure of objects and visual actions, allowing robots to predict future states and actions (Chen et al.[Link])

[1]

[2]

[3]

## Featured by a Geoscientist working on microscopic images (1)

**SET: Spectral Enhancement for Tiny Object Detection**

By Huixin Sun et al.
https://cvpr.thecvf.com/virtual/2025/poster/34394

Target: Detecting tiny objects

Proposal: HBS + API



Enhancing the distinctiveness of tiny features by suppressing high-frequency noise in the background

Adding adversarial perturbations to the feature map to enhance tiny object features during training.

**Achievements: +3.2% AP on AI-TOD (Tiny Object Detection) dataset**

## Featured by a Geoscientist working on microscopic images (2)

**Noise Calibration and Spatial-Frequency Interactive Network for STEM Image Enhancement**

By Hesong Li et al.    CVPR 2025 Open Access Repository

Target: Clearer observation by STEM*



Real Input | Theoretical GT | SrTiO₃ Structure

AtomSegNet Result Trained on GAN dataset | AtomSegNet Result Trained on TEMImageNet | AtomSegNet Result Trained on Our Dataset

*STEM: Scanning Transmission Electron **Microscope**

Proposal:



⇨ **Create Dataset**

**Achievements: More realistic visualization of minute structures**

LIMIT.LAB
https://limitlab.xyz/

[New Trends in Image Restoration and Enhancement workshop and associated challenges (NTIRE2025)](#)

NTIRE 2025 Challenge on Image Super-Resolution (x4)

❏ Participants should recover a high-resolution image from a single low-resolution input that is 4× smaller.

❏ Dataset
  ❏ The Challenge provides three official datasets: DIV2K, Flickr2K, and LSDIR.
  ❏ The LR-HR image pairs are generated using bicubic downsampling.
❏ Track
  ❏ The challenge evaluates performance on the DIV2K test set with two tracks:
  ❏ Restoration Track: Ranked by PSNR and SSIM on the Y channel.
  ❏ Perceptual Track: Ranked by a composite score from seven IQA metrics.

## [New Trends in Image Restoration and Enhancement workshop and associated challenges (NTIRE2025)](#)

NTIRE 2025 Challenge on Image Super-Resolution (x4)

[New Trends in Image Restoration and Enhancement workshop and associated challenges (NTIRE2025)](https://limitlab.xyz/)

NTIRE 2025 Challenge on Image Super-Resolution (x4)

## New Trends in Image Restoration and Enhancement workshop and associated challenges (NTIRE2025)

NTIRE 2025 Challenge on Image Super-Resolution (x4)

New Trends in Image Restoration and Enhancement workshop and associated challenges (NTIRE2025)

NTIRE 2025 Challenge on Image Super-Resolution (x4)

[New Trends in Image Restoration and Enhancement workshop and associated challenges (NTIRE2025)](#)

NTIRE 2025 Challenge on Image Super-Resolution (x4)

## Insights

❏ Transformer-based architectures remain a main-stream approach

❏ Integrating Mamba architectures for improved global context modeling

❏ Advanced training strategies, including multi-stage pipelines, progressive patch training, and CLIP-based semantic filtering, have boosted model generalization and robustness.

❏ Generative priors, particularly pre-trained diffusion models combined with CLIP-based perceptual losses, has achieved superior perceptual quality with minimal training.

## [ScanNet++ Novel View Synthesis and 3D Semantic Understanding Challenge](#)

This workshop focuses on the following key areas:

- ❏ High-fidelity, large-vocabulary 3D semantic scene understanding
- ❏ Novel view synthesis in large-scale 3D scenes

## ScanNet++ Novel View Synthesis and 3D Semantic Understanding Challenge

3D scene layout generation that combines LLMs with spatial reasoning

Enabling 3D scene generation technologies allows for applications in areas such as robotic simulation environments and game assets.

## ScanNet++ Novel View Synthesis and 3D Semantic Understanding Challenge

The trend is expected to shift toward 4D, which combines 3D and video.

## Workshop: Sight and Sound (paper session 1/2)

- CAV-MAE Sync: Improving Contrastive Audio-Visual Mask Autoencoders via Fine-Grained Alignment   main conference paper
  - Audio-visual masked autoencoder for audio-visual representation; temporal correspondences matter; improve CAV-MAE by a finer-grained temporal resolution; better audio-visual correspondence compared with ImageBind.

- Diagnosing and Treating Audio-Video Fake Detection
  - Deepfakes in audio-video; propose a benchmark dataset, a simple baseline, and evaluation protocols.

- UWAV: Uncertainty-weighted Weakly-supervised Audio-Visual Video Parsing main conference paper
  - Task: audio-visual video parsing; Proposed model: training the pseudo label models with larger training set; training with uncertainty-weighted feature mixup; Evaluation dataset: LLP dataset which contains both audio events, visual events, and audio-visual events.

- STM2DVG: Synthetically Trained Music to Dance Video Generation leveraging Latent Diffusion Framework
  - Pose-conditioned LDM; Synthetic dataset generation pipeline; music-to-pose encoder.

- Seeing Speech and Sound: Distinguishing and Locating Audio Sources in Visual Scenes main conference paper
  - Task: simultaneous audio-visual grounding, including overlapping spoken language and non-speech sounds; Method:

LIMIT.LAB

https://limitlab.xyz/

## Workshop: Sight and Sound (paper session 2/2)

- ❏ [AVS-Net: Audio-Visual Scale Net for Self-supervised Monocular Metric Depth Estimation](#)
  - ❏ Concept: using echoes for enhancing depth estimation; Especially using echos to obtain higher scale accuracy.

- ❏ [SAVGBench: Benchmarking Spatially Aligned Audio-Video Generation](#)
  - ❏ Concept: spatially aligned audio-video generation; Propose dataset (with ambisonics audio, 360-degree videos), with human speeches and instrument sounds; Method: MM Diffusion based.

- ❏ [BGM2Pose: Active 3D Human Pose Estimation with Non-Stationary Sounds](#)
  - ❏ Concept: using non-stationary sound to aid active 3D human pose estimation;

- ❏ [Visual Sound Source Localization: Assessing Performance with Both Positive and Negative Audio](#)
  - ❏ Visual sound source localization; Motivation: enhance robustness to offscreen and white noise; Extended datasets by adding negative audio; Introduced several evaluation metrics;

- ❏ [VGGSounder: Audio-Visual Evaluations for Foundation Models](#)
  - ❏ Introduce a comprehensive av benchmark, which is human labeled and each video has annotated classes in each modality.

## Workshop: Sight and Sound (invited speaker)

- ❏ Learning to Infer Audio-Visual Attention in Social Communication (James Rehg)
    - ❏ Egocentric perception is important in social communication
        - ❏ Paper: Listen to Look into the Future: Audio-Visual Egocentric Gaze Anticipation (ECCV2024)
            - ❏ Task: gaze anticipation
            - ❏ Dataset (Ego4D Social, Aria)
            - ❏ Method: using video and audio, spatial and temporal fusion
    - ❏ Vision foundation model for gaze estimation
        - ❏ Paper: Gaze-LLE: Gaze Target Estimation via Large-Scale Learned Encoders (CVPR 2025)
            - ❏ Task: recognition human gazes in all people in videos
            - ❏ Method: utilizing foundation model + finetuning decoder part
    - ❏ Group conversation recognition
        - ❏ Paper: The Audio-Visual Conversational Graph: From an Egocentric-Exocentric Perspective (CVPR2024)
            - ❏ Motivation: ego-centric, exo-centric both important
            - ❏ Task: ego-exocentric conversational graph prediction
            - ❏ Model: multiple encoders (image, audio), plus cross, and self-attention
    - ❏ Social gestures, and other social behaviors, combined with computer vision tasks
        - ❏ Paper: SocialGesture: Delving into Multi-person Gesture Understanding (CVPR2025)
            - ❏ A benchmark to recognize and localize gestures
            - ❏ Current multimodal LLMs, foundation models still fall short in recognizing multiple person social interactions
        - ❏ Paper: Werewolf Among Us: A Multimodal Dataset for Modeling Persuasion Behaviors in Social Deduction Games
            - ❏ Persuasion strategy prediction task
        - ❏ Other datasets: speaking target identification,

## Workshop: Sight and Sound (invited speaker)

- ❏ Sight and Sound with Large Language Models: Applications to Video Dubbing and Spatial Sound Understanding ([David Harwath](#))

- ❏ Expanding the ability of multimodal LLMs in two applications

- ❏ Video dubbing

  - ❏ Paper: VOICECRAFT: Zero-Shot Speech Editing and Text-to-Speech in the Wild

    - ❏ Task: voice-cloning TTS

    - ❏ Issues: not enough talking-head datasets

    - ❏ Approach: firstly using language-based data to train text-to-speech backbone; the second, using video data to finetune

    - ❏ Other ideas: tokenizing speech with neural codecs

  - ❏ Paper: VoiceCraft-Dub: Automated Video Dubbing with Neural Codec Language Models

    - ❏ Task: dubbing, which should generate audio aligned to person, text, and video

    - ❏ Method: initialize from VoiceCraft model; add lip/face encoders and finetune on video-based datasets

- ❏ Spatial sound understanding

  - ❏ Paper: BAT:Learning to Reason about Spatial Sounds with Large Language Models

    - ❏ Task: spatial sound recognition; sound event detection; direction; distances, …

    - ❏ Dataset: SoundSpaces 2.0 simulator

    - ❏ Method: adding spatial sound encoding to LLM + pretraining (Spatial-AST model)

## Workshop: Sight and Sound (invited speaker)

- ❏ Learning Sight and Sound through Generative Models (Ziyang Chen)
  - ❏ Paper: Diff-Foley: Synchronized Video-to-Audio Synthesis with Latent Diffusion Models
    - ❏ LDF for sound generation from video;
  - ❏ Paper: Video-Guided Foley Sound Generation with Multimodal Controls (CVPR 2025)
    - ❏ multimodal conditional; joint training using paired video+text+audio / text+audio data.
    - ❏ Very impressive generation result
  - ❏ Paper: Composing Images and Sounds on a Single Canvas (NeurIPS 2025)
    - ❏ Intersection between image representation and spectrogram
  - ❏ Idea: images that sound (image looking like the image caption, and sound like the sound caption)
  - ❏ Method: using the diffusion model

## Workshop: 3D Vision Language Models (VLMs) for Robotic Manipulation: Opportunities and Challenges (invited speaker)

- ❏ Building vision-language maps for embodied AI (Angel Chang)
  - ❏ Semantic Mapping in Indoor Embodied AI (survey paper)



- ❏ Building maps for object-centric navigation
  - ❏ Multi-Object Navigation (MultiON) Task, LangNavBench, GOAT-Bench
  - ❏ Multi-Layered Feature Map
    - ❏ Two phase: explore-and-explore to get object information (build map) -> explore only (find goal in map)
  - ❏ Paper: Zero-shot Object-Centric Instruction Following: Integrating Foundation Models with Traditional Navigation
    - ❏ Using foundation models; Proposed Language-Inferred Factor Graph for Instruction Following, a graph structure, leading to more structured navigation and good accuracy.

**Workshop: 3D Vision Language Models (VLMs) for Robotic Manipulation: Opportunities and Challenges (invited speaker)**

- ❏ Hierarchical Action Models for Open-World 3D Policies (<u>Dieter Fox</u>)
  - ❏ Perceiver Actor
    - ❏ Anticipating the 3D poses / positions before executing actions
    - ❏ Paper: RVT: Robotic View Transformer for 3D Object Manipulation
      - ❏ Building explicit 3D representation by generating robotic observation by combining multiple viewpoint information
    - ❏ Paper: RVT2: RVT-2: Learning Precise Manipulation from Few Demonstrations
      - ❏ Predicting the next key-frame pose
  - ❏ Guiding 3D policies by utilizing vision language models
    - ❏ Input -> Vision-Language-Action Model -> Context for low-level policy -> control
    - ❏ Paper: HAMSTER: Hierarchical Action Models For Open-World Robot Manipulation (ICLR 2025)
      - ❏ VLM generates 2D trajectory sketch  for policy learning
      - ❏ 3D policy models for 3D policy training
      - ❏ Shows promising generation ability in real-world evaluation of generalization

## Workshop: 3D Vision Language Models (VLMs) for Robotic Manipulation: Opportunities and Challenges（invited speaker）

- ❏ Genesis: An Unified and Generative Physics Simulation for Robotics (Chuang Gan)
  - ❏ How far are we from an embodied generalist agent?
    - ❏ Generalist agents: multi-modal; multi-task; multi-environment.
    - ❏ At least, we need a world model
  - ❏ How to model physical world
    - ❏ Generative AI + Differentiable Physical engines
    - ❏ Forward Simulations Using Differentiable Physics
    - ❏ Joint geometry and physics estimation from multi-view videos
    - ❏ Forward simulation with actions
    - ❏ Result is better than RL-based learning
  - ❏ Applications
    - ❏ Soft-object animation; Thin-shell manipulation
  - ❏ Can we scale up the above approach?
    - ❏ Genesis: general-purpose physics simulator; Fully differentiable;
      - ❏ To scale up data generation, the authors proposed RoboGen, capable of task proposal (using language, using language as prompt), scene generation, training supervision generation.
      - ❏ How to solve the issue that human tools could not be suitable for robotics? -> Generative tools for robotics

## Workshop: 3D Vision Language Models (VLMs) for Robotic Manipulation: Opportunities and Challenges（spotlight papers）

- ❏ The One RING: A Robotic Indoor Navigation Generalist
  - ❏ Universal navigation policy for all different robotics? – different robots observe and behave differently. The authors proposed One Ring, both data and model; Method: large-scale training + RL fine-tuning on random embodiments; Trained on sim but can be adapted to real-env.
- ❏ Manual2Skill: Learning to Read Manuals and Acquire Robotic Skills for Furniture Assembly Using Vision-Language Models
  - ❏ Manipulation skills learning by utilizing vision-language models; Use 2D images and language as input, generate hierarchical graph for assembly.
- ❏ Agentic Language-Grounded Adaptive Robotic Assembly
  - ❏ Issue: the ability to adapt to different model shape and structure change accordingly is essential in assembly applications; Resolve: using LLM to generate hierarchical guide for adapting assembly skills for different parts and models change.
- ❏ ZeroMimic: Distilling Robotic Manipulation Skills from Web Videos
  - ❏ Issue: datasets for robotic training; Resolve: learn general skill policies from online videos; Model: step1: knowing what human hands are doing; step2: training BC policies to operate human hands; step3: deploying human arm policies on robot arms. Foundation models facilitate distilling human videos into robot policies.

## Workshop: Visual Generative Modeling: What's After Diffusion?

- ❏ After Diffusion Models (Bill Freeman)
  - ❏ Problems of diffusion models
    - ❏ Slow, not easy control, not easy editing, …
  - ❏ Possible methods that might work
    - ❏ Conventional Graphics, hiring human artists, Simple explainable generative algorithms, Engineering, Better reasoning models, …
  - ❏ Conventional Graphics
    - ❏ Why not: slow, complicated, lots of manual work
  - ❏ Simple explainable generative algorithms
    - ❏ Paper: Infinite Images: Creating and Exploring a Large Photorealistic Virtual Space
      - ❏ Unwrap the motion into an infinite image
    - ❏ Paper: WonderWorld: Interactive 3D Scene Generation from a Single Image (CVPR 2025)
  - ❏ Computer vision by asking nicely
    - ❏ Foundation model that has enough knowledge about the physical world
    - ❏ Paper: Alchemist: Parametric Control of Material Properties with Diffusion Models (CVPR 2024)
    - ❏ Coaxing LLM to do what you ask, infuriating but sometimes powerful
  - ❏ Engineering better diffusion models
    - ❏ Paper: Improved Distribution Matching Distillation for Fast Image Synthesis (NeurIPS 2024)
    - ❏ Distillation of diffusion model

## Workshop: Visual Generative Modeling: What's After Diffusion?

- ❏ Controllable, Intuitive Generation of 4D Objects and Scenes (1/2) ([Jiajun Wu](#))
    - ❏ Perception and Generation beyond 2D pixels
        - ❏ De-render to physical object intrinsics – allow better controllableness
    - ❏ Idea1: intrinsics as inductive bias for SSL
        - ❏ De-Render to intrinsics and then re-render to images
        - ❏ Physical intrinsics (lighting, 3D shape, camera, ..)
        - ❏ Paper: Seeing a Rose in Five Thousand Ways (CVPR 2023)
            - ❏ Single image, multiple object instances (like bouquet) –all roses have similar intrinsics
            - ❏ Make it possible to controllability
        - ❏ Paper: Learning the 3D Fauna of the Web (CVPR 2024)
            - ❏ From a single category to multiple categories, extending the above
    - ❏ Idea2: intrinsics as distilling targets for FMs (only vision)
        - ❏ Paper: Birth and Death of a Rose (CVPR 2025)
            - ❏ Generating temporal object intrinsics—temporally evolving sequences of object geometry, reflectance, and texture, such as a blooming rose—from pre-trained 2D foundation models.
        - ❏ Paper: PhycsDream: From Object Motion to Action (ECCV 2024)
            - ❏ Action-conditioned dynamics prediction of 3D objects

## Workshop: Visual Generative Modeling: What's After Diffusion?

- ❏ Controllable, Intuitive Generation of 4D Objects and Scenes (2/2) ([Jiajun Wu](#))
  - ❏ Idea2A: intrinsics as distilling targets for FMs (vision and language)
    - ❏ From vision-language models observations to reconstruction
    - ❏ Paper: WonderJourney: Going from Anywhere to Everywhere (CVPR 2024)
      - ❏ Model illustration: using LLMs to generate long sequence of scene descriptions, a text-driven point cloud generation pipeline to synthesize 3D scenes, and a VLM to verify generated examples
    - ❏ Paper: WonderWorld: Real-Time, Interactive 3D Scene Generation (CVPR2025)
      - ❏ User can specify where and what to generate in an interactive manner
    - ❏ Paper: WonderPlay: Dynamic 3D Scene Generation from a Single Image and Actions
      - ❏ Allow editing and manipulation in 4D manner
    - ❏ Paper: The Scene Language: Representing Scenes with Programs, Words, and Embeddings (CVPR 2025)
      - ❏ Chessboard generation as an example, allowing generating scenes expressively, allowing high-level interactive editing
      - ❏ Structure: program function dependency
      - ❏ Semantics: natural language words
      - ❏ Visual identity: Neural embeddings
      - ❏ Allowing Text-to-3D/4D Generation and Editing
- ❏ What's the next:
  - ❏ Questions: what needs to be modeled? How do they help visual generation? What's the roles of FMs
  - ❏ Rendering/conditioning based on causal, physical, and universal object intrinsics
  - ❏ Effective use of the in-the-wild data, enhanced interpretability

## Workshop: Visual Generative Modeling: What's After Diffusion?

- ❏ Language as a Visual Format ([Phillip Isola](#))
  - ❏ How similar is the way a vision model sees the world to the way a language model sees the world
  - ❏ Paper: The Platonic Representation Hypothesis (ICML, 2024)
    - ❏ Characterizing representations using kernels
    - ❏ Comparing sim (DINO, Llama)
    - ❏ The bigger the DINO model is, the similarity is higher
    - ❏ More detailed captions align better with vision representations
  - ❏ Paper: Cycle Consistency as Reward: Learning Image-Text Alignment without Human Preferences
    - ❏ Image and language cycle consistency to improve image-text alignment
    - ❏ Light-weight, fast, differentiable
    - ❏ Dataset: Cycle consistency preference collection (CyclePrefDB)
      - ❏ Useful for evaluating and improving image-text alignment
      - ❏ Direct Preference Optimization using CyclePrefDB-T2I
      - ❏ Usage: pip install cyclereward
- ❏ What's after diffusion
  - ❏ Language as an alternative to pixels but requires visually descriptive language

## Workshop: VLMs-4-All 2025

❏ Re-scaling cultural knowledge in a UK heritage and museum context (Maya Indira Ganesh)
  ❏ AI for cultural heritage hub (ArCH)
    ❏ Important question: how to use current AI to let people be more engaged in cultural heritage recognition
  ❏ GLAM sector and cultural heritage collection challenges
    ❏ Challenges: Accessibility, reconstructing fragmentary and dispersed objects, requiring expert knowledge, legacy technical systems across multiple providers, legacies of colonialism and contested provenance of artifacts and cultural knowledge
    ❏ Some other issues: negative perception of AI, uncertainty, errors, digitization, undigitized backlog, legacy infrastructure, provisional semantics
  ❏ Street observatories of everyday AI
    ❏ Street recognition and embodied tasks in streets still have lots of unresolved issues

## Workshop: VLMs-4-All 2025

❏ Richer Outputs for Richer Countries: Geographical Disparities in Language and Image Generation (1/2) ([Danish Pruthi](#))
  ❏ Mononyms (single word names)
    ❏ Common in some modern societies (like India, Myanmar)
    ❏ Digital forms rarely account for single names
    ❏ Issues: might not be recognized by recent LLMs
    ❏ Other AI system issues: biases related to gender, race, geographical
  ❏ Generating Representative
    ❏ Generating an image / images representative of visual concepts, which is important in various applications and real-world scenarios
  ❏ Issue: Are current models generating geographical representative?
    ❏ Paper: Inspecting the Geographical Representativeness of Images from Text-to-Image Models (ICCV 2023)
      ❏ Current models show low-level of geographical representativeness



  ❏ Paper: Easily Accessible Text-to-Image Generation Amplifies Demographic Stereotypes at Large Scale (FAccT 2023)
    ❏ A broad range of ordinary prompts produce stereotypes, including prompts simply mentioning traits, descriptors, occupations, or objects
    ❏ Stereotypes are present regardless of whether prompts explicitly mention identity and demographic language or avoid such language.

*LIMIT.LAB*

## Workshop: VLMs-4-All 2025

❏ Richer Outputs for Richer Countries: Geographical Disparities in Language and Image Generation (2/2) ([Danish Pruthi](#))
   ❏ Issue: Are current models generating geographical representative? - No
      ❏ Paper: Where Do Images Come From? Analyzing Captions to Geographically Profile Datasets (left image)
         ❏ Task: Geographically profiling datasets, given an image-caption pair, map them to a location
         ❏ Geo-profiling LAION2B-EN proves heavy geographical biases, data distribution greatly biased where Wealthy countries have more data



   ❏ Language oppression: erasure of minority languages
      ❏ Paper: Geographical Erasure in Language generation (EMNLP 2023)
         ❏ Many countries suffer from geographical erasure
      ❏ Paper: WorldBench: Quantifying Geographic Disparities in LLM Factual Recall (right image)
         ❏ Assess the ability of large language models (LLMs) to recall factual information about specific countries
   ❏ Paper: Richer Output for Richer Countries: Uncovering Geographical Disparities in Generated Stories and Travel Recommendations
      ❏ Reveal the issues of geographical erasure, regarding to applications such as travel recommendations and story generation

## Workshop: VLMs-4-All 2025

❏ Building Culturally Aware Multilingual LMM Benchmarks (1/2) ([Fahad Shahbaz Khan](#))
  ❏ Efforts 1:  multilingual multimodal understanding benchmarks
    ❏ Paper: All Languages Matter: Evaluating LMMs on Culturally Diverse 100 Languages (CVPR2025) (left image)
      ❏ Culturally Aware Multilingual Image LMM Benchmark
      ❏ 100 languages, including many low-resource ones, native-language experts (800 hours annotations)
      ❏ 13 cultural aspects covered in the benchmark, such as heritage, customs, architecture, literature, music, and sports.
      ❏ Generation pipeline with AI generated QAs with human validation
    ❏ Paper: A culturally-diverse multilingual multimodal video benchmark and model (right image)
      ❏ Issue: Most video datasets are concentrating in English, models struggle with low-resource languages
      ❏ Propose: Culturally Aware Multilingual Video LMM Benchmark

## Workshop: VLMs-4-All 2025

❏ Building Culturally Aware Multilingual LMM Benchmarks (2/2) (Fahad Shahbaz Khan)
  ❏ Efforts 2: Arabic LMM benchmarks
    ❏ Paper: CAMEL-Bench: A Comprehensive Arabic LMM Benchmark (left image)
      ❏ CAMEL-Bench, comprises eight diverse domains and 38 sub-domains including, multi-image understanding, complex visual perception, handwritten document understanding, video understanding, medical imaging, plant diseases, and remote sensing-based land use understanding
      ❏ Arabic language has over 400M users, however, even the closed-source GPT-4o achieving an overall score of 62%
    ❏ Paper: ARB: A Comprehensive Arabic Multimodal Reasoning Benchmark (right image)
      ❏ The first benchmark designed to evaluate step-by-step reasoning in Arabic across both textual and visual modalities.
      ❏ ARB spans 11 diverse domains, including visual reasoning, document understanding, OCR, scientific analysis, and cultural interpretation.

## Workshop: VLMs-4-All 2025

- ❏ Challenges with Geo-Cultural Understanding and Generation ([Roopal Garg](#))
  - ❏ GeoCultural Multimodal Content Understanding and Generation
  - ❏ Issue: current models fall short on a truly global user base
    - ❏ Impact: affects accessibility, fairness, utility, and trust of billions of users worldwide
  - ❏ Data quality gap – comprehensiveness vs. cultural specificity
    - ❏ Not having enough high-quality scaled datasets
    - ❏ Data generation / collection: collect data that highlights cultural distinctiveness, e.g., food, festivals.
    - ❏ Generate Locale Aware human annotations, instead of direct translations, using locale reference
    - ❏ Search from the internet to get correct labels, and use local person annotations or validations
  - ❏ Generation model for enhanced geo-cultural understanding
    - ❏ Issue: image generation could be culture-related, text-to-image generation from english translations contain multiple issues and is less accurate for concepts with different cultural backgrounds
    - ❏ Model structure: prompt expansion (get local language and modification based on the context of the target language) before generation

## Workshop: VLMs-4-All 2025

- ❏ Building geo-diverse model ([Olga Russakovsky](#))
  - ❏ Current issue: Lack of cultural diversity in datasets
  - ❏ GeoDE: Geographically Diverse Evaluation Dataset (NeurIPS2023) (left image)
    - ❏ Partnered with Appen to solicit photographs from people around the world
    - ❏ 61940 images from 6 different regions
    - ❏ Generated using crowdsourcing where each image has consent and no recognizable people
    - ❏ Findings:
      - ❏ Crowd-sourced images look different in different regions (right image)
      - ❏ Features representations between crowd-sourced and web-scraped images are different
      - ❏ GeoDE can find gaps in models, for example, CLIP shows stereotypes and errors in recognizing non-western objects, building
      - ❏ Training on GeoDE improves performance

## ExpertAF: Expert Actionable Feedback from Video

- **Summary:** This work proposes a task that generates video commentary, movement guidance (text), and correct posture (skeleton) from amateur videos (e.g., playing soccer) and corresponding pose (skeleton) data. From the Ego-Exo4D dataset for skill learning, they paired similar samples (amateur video + pose, expert video + pose) and used an LLM to generate movement guidance, thereby semi-automatically constructing a dataset. The method is simple, consisting of separate modules: an encoder for each input, an LLM-based motion guidance generator, and a Retrieval-Augmented Generator-based pose generator.
- **Novelty:** Task definition and dataset.
- **Thoughts:** The part that also outputs poses is interesting. There are still many challenges in tasks that aim to generate high-quality and visually meaningful feedback. Both the method and the dataset seem to have room for improvement. It's also worth exploring whether training with a large amount of data can produce precise feedback. Generating demonstrations for group activities may be particularly challenging.

## FICTION: 4D Future Interaction Prediction from Video

- ❏ **Summary:** This paper proposes *FICTION*, a method for 4D estimation from video—i.e., determining *when*, *where*, *which* object is being interacted with, and *how* (its pose). FICTION uses two VAE structures to estimate location and pose respectively from the embedding of past observations (see figure on the right).
- ❏ **Novelty:** While previous studies have explored estimating Video Scene Graphs from video, this work extends that task into 4D estimation.
- ❏ **Thoughts:** Research combining video and 3D is increasing. Papers from Kristen Grauman's group often use simple and easy-to-understand methods, especially for task-setting works—where they tend to first establish baselines using straightforward approaches. Also, since the same group is working on other 4D generation tasks, it seems they are shifting from video-only to video + 3D (i.e., 4D) representations.

## Video-3D LLM: Learning Position-Aware Video Representation for 3D Scene Understanding

- ❏ **Summary:** This work proposes *Video-3D LLM*, a novel MLLM (Multimodal Large Language Model) method capable of 3D recognition. It treats 3D as dynamic video and aligns video representations with 3D representations. The model achieves high accuracy on several 3D recognition benchmarks.
- ❏ **Novelty:** The key novelty lies in recognizing 3D information from video, and enabling MLLMs to understand 3D using this approach.
- ❏ **Thoughts:** The ability to perceive 3D from video alone feels highly practical. In real-world applications, it's not always necessary to have explicit 3D input if the system can understand 3D from videos. Being able to recognize 3D using only video makes the method much more user-friendly and applicable.



LIMIT.LAB

## METASCENES: Towards Automated Replica Creation for Real-world 3D Scans

- ❏ **Summary:** This work proposes a new dataset, *METASCENES*, and a retrieval-based method called *Scan2Sim*. METASCENES is built upon ScanNet, augmented with human annotations and enhanced through physical optimization by replacing ScanNet assets with CAD models. Experiments demonstrate the effectiveness of METASCENES in improving accuracy.
- ❏ **Novelty:** The work introduces a pipeline that easily converts real scanned data into simulation data. This significantly reduces the cost of generating 3D indoor environment datasets.
- ❏ **Thoughts:** It would be interesting to see comparisons with other methods for automatic creation of indoor 3D simulation data. Text-based generation methods might be more user-friendly in some cases.

## ARKit LabelMaker: A New Scale for Indoor 3D Scene Understanding

- **Summary:** This work proposes *ARKit LabelMaker*, a large-scale 3D semantic segmentation dataset—over three times larger than existing datasets. It also introduces an automated pipeline to generate ARKit LabelMaker, leveraging several high-performing models (OVSeg, Grounded-SAM, InternImage, Mask3D). Methods pre-trained on ARKit LabelMaker outperform those trained on existing datasets.
- **Novelty:** The dataset's scale is a key contribution, along with the automated pipeline that enables further scaling of data generation.
- **Thoughts:** If this approach could also scale language + 3D data together, performance might improve even further.

## MARVEL-40M+: Multi-Level Visual Elaboration for High-Fidelity Text-to-3D Content Creation Creation

- ❏ **Summary:** This work proposes *MARVEL-40M+*, a large-scale dataset of text and 3D models (40 million text entries and 8.9 million 3D shapes). It introduces an automated dataset generation pipeline using a multimodal LLM to generate tags and detailed descriptions for each 3D model. The authors also propose a two-stage *Text-to-3D* generation model: Stage 1: Text-to-Image, and Stage 2: Image-to-3D.
- ❏ **Novelty:** The key novelty lies in performing both detailed 3D annotation and 3D model generation. The strong experimental results further validate the approach.
- ❏ **Thoughts:** It would be interesting to see tasks that combine 3D prompts (bounding boxes, segmentation, text) with 3D recognition.

## IDEA-BENCH: HOW FAR ARE GENERATIVE MODELS FROM PROFESSIONAL DESIGNING?

- ❏ **Summary:** A new benchmark called IDEA-BENCH was created to examine how closely current generative models can perform design work in a "professional designer" setting. The benchmark contains 100 tasks spanning text-to-image, image-to-image, and image-editing scenarios, and it was used to evaluate several prominent models. The study offers a detailed analysis of the shortcomings those models exhibit when asked to solve genuine design problems.

- ❏ **Novelty:** IDEA-BENCH assesses real-world design tasks along multiple axes. Even powerful systems such as DALL-E 3 and InstructPix2Pix achieve only about 20 points out of 100 on this benchmark, revealing systematic weaknesses that had not been fully documented.

- ❏ **Thoughts:** Evaluation tasks that gauge a model's proximity to human expertise are becoming more common. If a dataset can reliably expose crucial failure modes, then a moderate dataset size could be enough. Developing robust evaluation schemes for generative models remains an intriguing research direction.

## Joint Vision-Language Social Bias Removal for CLIP

❏ **Summary:** A method is introduced for removing gender, race, and age bias from CLIP embeddings. Previous approaches that strip biased attributes from embeddings often degrade downstream task performance. Through a comprehensive analysis, the study presents a technique that debiases both the image and text branches of CLIP, achieving strong bias mitigation while largely preserving accuracy on downstream tasks.

❏ **Novelty:** Debiasing CLIP remains relatively under-explored despite its importance. The work offers a thorough critique of existing methods and addresses their limitations, delivering an approach that balances fairness and task performance.

❏ **Thoughts:** Investigating bias in large (multi)-modal language models is a pressing challenge, and this contribution highlights both the difficulties and the potential of targeted debiasing strategies.

## HD-EPIC: A Highly-Detailed Egocentric Video Dataset

- **Summary:** The study presents HD-EPIC, a fine-grained video dataset enriched with recipe steps, atomic actions, ingredient lists with nutritional values, object-tracking labels, audio annotations, and fully aligned 3-D digital-twin data. Even the strong Gemini Pro model achieves only 37.6 % accuracy on the accompanying benchmark, underscoring the dataset's difficulty.
- **Novelty:** HD-EPIC extends existing resources by combining fine-grained video annotations with explicit 3-D twins, creating a unified test bed for vision-language and embodied-AI research. The tight linkage between 2-D video and 3-D geometry is a particularly novel contribution.
- **Thoughts:** The results highlight how current multi-modal language models still fall short on detailed video understanding. Producing fine-grained, 3-D-aligned labels is labor-intensive, yet the resulting dataset should prove valuable for a wide range of video, robotics, and embodied-reasoning tasks.

## HarmonySet: A Comprehensive Dataset for Understanding Video-Music Semantic Alignment and Temporal Synchronization

- ❏ **Summary:** HarmonySet is introduced as an instruction-tuning dataset for training and evaluating multi-modal language models on joint music–video understanding. The collection supports assessment of rhythmic synchronization, emotional alignment, thematic coherence, and cultural relevance. All annotations are generated automatically with the aid of existing MLLMs, enabling large-scale coverage.
- ❏ **Novelty:** Very few benchmarks explicitly target the fusion of musical and visual semantics for MLLMs. HarmonySet fills this gap and offers new evaluation axes that go beyond conventional audiovisual alignment.
- ❏ **Thoughts:** There is a visible trend toward building datasets and metrics for challenging domains such as music and design, where subjective judgments make evaluation difficult. HarmonySet represents a timely step in that direction.

## ComfyBench: Benchmarking LLM-based Agents in ComfyUI for Autonomously Designing Collaborative AI Systems

- ❏ **Summary**: A benchmark called ComfyBench and an accompanying method named ComfyAgent are introduced for automatically generating—given a high-level task—the required tools, step-by-step workflow, and intermediate prompts. The resulting "flow" is expressed in an easily editable code format, and ComfyAgent reaches accuracy comparable to the O1-Preview system.
- ❏ **Novelty:** The work extends visual-programming ideas to a newer tool ecosystem; although the conceptual leap is modest, matching O1-Preview's performance is noteworthy.
- ❏ **Thoughts:** Building agents by combining predefined toolsets with LLM-generated prompts reflects an interesting trade-off and resembles modern visual-programming paradigms.

## MicroVQA: A Multimodal Reasoning Benchmark for Microscopy-Based Scientific Research

- ❏ **Summary:** The paper presents MicroVQA, a microscopy-based visual-question-answering dataset. Unlike conventional VQA corpora, MicroVQA is organized around the full scientific workflow—observation, hypothesis generation, and experimental verification—and was curated by professional biologists. Several multi-modal language models are evaluated on the benchmark, revealing a clear advantage for larger models.
- ❏ **Novelty:** MicroVQA goes beyond simple question–answer pairs by embedding scientific reasoning steps into the data structure. Evaluation shows that most errors stem from incorrect visual recognition, indicating that interpreting microscopy images remains a bottleneck.
- ❏ **Thoughts:** "AI for Science" resources are on the rise, and datasets that mirror the entire research process, such as MicroVQA, appear especially valuable. Results suggest that domain-specific imagery (microscopy, remote sensing, etc.) still suffers from insufficient pre-training in current models.



Table 1. MicroVQA benchmark attributes.

| Dataset feature | Value |
|---|---|
| Total questions | 1,042 |
| Multi-image questions | 423 |
| Avg. MCQ question length | 66 |
| Avg. MCQ answer length | 15 |
| Avg. raw question length | 158 |
| Avg. raw answer length | 52 |
| Unique image sets | 255 |
| Image Modalities | Light, Fluoro, Electron |
| Image Scales | Tissue, Cell, Subcell, Atomic |
| Organisms | 31 |
| Research areas | 33 |
| Expert question creators | 12 |
| Time to create 1 question | 30-40 mins |
| Time to quality check 1 MCQ | 5 mins |

LIMIT.LAB

## Automated Generation of Challenging Multiple-Choice Questions for Vision Language Model Evaluation

- ❏ **Summary:** The paper introduces AutoConverter, an MLLM-based method that automatically transforms existing open-ended VQA datasets into multiple-choice format.  The converted data enable the creation of much larger multiple-choice VQA benchmarks, making it easier to evaluate multi-modal language models at scale.

- ❏ **Novelty:** AutoConverter's main contribution lies in repurposing open-ended VQA data for multiple-choice evaluation.  While useful, the conceptual advance is relatively modest.

- ❏ **Thoughts:** The same idea could plausibly be extended to video VQA.  Although image and video datasets are abundant, comprehensive benchmarks for 3-D recognition remain scarce.  In the long run, it would be valuable to develop AI systems that can identify their own knowledge gaps, acquire relevant data, and continue learning autonomously.



(a) *AutoConverter* framework.

## Unbiased General Annotated Dataset Generation

❏ **Summary:** A data-generation method called ubGen is introduced to improve the generalization of image-recognition models. The approach generates images along directions in CLIP's semantic space and employs a quality–text-alignment function to ensure both fidelity and diversity. Models trained on ubGen-augmented data achieve markedly better transfer performance across several benchmark datasets.

❏ **Novelty:** ubGen can produce highly transferable training images using nothing more than the category names, offering a lightweight alternative to large-scale manual collection.

❏ **Thoughts:** Work on boosting diffusion-model expressiveness, often by leveraging MLLMs, appears to be gaining momentum, and ubGen fits neatly into this emerging trend.

## FaceBench: A Multi-View Multi-Level Facial Attribute VQA Dataset for Benchmarking Face Perception MLLMs

❏ **Summary:** A dataset named FaceBench is introduced for multi-view, fine-grained face recognition. It contains 210 annotated facial attributes and more than 70 000 question, answer pairs. A model pre-trained on FaceBench attains accuracy comparable to GPT-4o and Gemini on the benchmark.

❏ **Novelty:** FaceBench provides substantially finer facial-attribute coverage than previous resources, enabling more detailed recognition and analysis.

❏ **Thoughts:** While FaceBench targets static face images, a complementary video corpus for nuanced expression recognition would also be valuable. Growing demand is evident for fine-grained datasets across many domains, though it is somewhat surprising that the FaceBench-trained model does not surpass GPT-4o.

LIMIT.LAB

https://limitlab.xyz/

## UnCommon Objects in 3D

- **Summary:** A large-scale 3-D object dataset named uCO3D is introduced. Spanning more than 1 000 object categories, each sample includes accurate structure-from-motion camera poses, a point cloud, and a 3-D Gaussian-splat reconstruction. Experiments show that models trained on uCO3D achieve markedly better performance in novel-view synthesis and 3-D reconstruction than when trained on MVImageNet or CO3Dv2.
- **Novelty:** uCO3D surpasses existing datasets in both scale and annotation quality, offering denser geometry and more object diversity than previous resources.
- **Thoughts:** Pairing such high-fidelity geometry with strong language supervision could be an interesting next step, although complex shapes and textures remain difficult to describe textually.

## HumanDreamer: Generating Controllable Human-Motion Videos via Decoupled Generation

- ❏ **Summary:** This paper introduces a diffusion-based video generation method that produces both poses and corresponding videos from text inputs in a two-stage process: first generating poses from text, then generating video from the predicted poses. To enable the text-to-pose stage, a large-scale, semi-automatically constructed dataset is presented. Additionally, a new loss function, LAMA loss, is proposed for the text-to-pose model to better align predicted poses with pose features.
- ❏ **Novelty:** Key contributions include the construction process and release of a large-scale Text-to-Pose dataset, as well as the introduction of the LAMA loss tailored for pose prediction from text.
- ❏ **Thoughts:** Text-to-pose generation has broad potential applications. This work also follows a growing trend of using pre-existing models and large language models (LLMs) to automate dataset construction.

## Apollo: An Exploration of Video Understanding in Large Multi-Modal Models

❏ **Summary:** This Meta research conducts a comprehensive experimental study and comparison of Video LLM design choices, culminating in the introduction of the Apollo model series. Apollo achieves higher accuracy compared to models with similar computational cost.

❏ **Novelty:** Key findings include the observation of scaling consistency—performance trends seen on relatively large datasets continue to hold on even larger ones. Additionally, FPS-based frame sampling is shown to outperform uniform sampling, yielding performance gains proportional to the number of frames, which is not the case for uniform sampling. At the time of publication, models using InternVideo2 + SigLIP SO400M features achieved the best results.

❏ **Thoughts:** The study involves an extensive amount of experimentation and provides valuable insights into the design of video recognition models.

## ByTheWay: Boost Your Text-to-Video Generation Model to Higher Quality in a Training-free Way

- ❏ **Summary:** This work introduces a training-free method to improve generation consistency and motion magnitude in diffusion-based video generation. Through an analysis of temporal attention, the authors identify that differences between temporal attention maps strongly affect temporal consistency. By adjusting these maps to enhance consistency, overall generation stability improves. Additionally, the energy of temporal attention maps is found to influence motion magnitude, leading to a method that amplifies motion strength. The proposed technique significantly improves the output quality of existing models such as AnimationDiff and VideoCrafter2.
- ❏ **Novelty:** The study offers a comprehensive analysis of temporal attention maps and introduces a novel, training-free method for enhancing consistency and motion dynamics in generated videos.
- ❏ **Thoughts:** It would be valuable to verify whether the proposed method generalizes to other diffusion-based video generation models.

## Holmes-VAU: Towards Long-term Video Anomaly Understanding at Any Granularity

- ❑ **Summary:** This paper introduces **Holmes-VAU**, a large-scale dataset designed to recognize anomalies in videos across **multiple temporal granularities**. The dataset is semi-automatically constructed using existing multi-modal language models (MLLMs). In addition, a method combining a **temporal sampler** with an MLLM-based model is proposed for anomaly detection at various time scales.

- ❑ **Novelty:** The key contribution lies in formulating anomaly recognition across **different temporal granularities**, which offers a new perspective compared to standard fixed-scale video analysis.

- ❑ **Thoughts:** The idea of analyzing video at varying temporal resolutions is clearly structured and meaningful. However, since related datasets like temporal localization benchmarks already exist, the degree of novelty may be somewhat limited.

## LongVALE: Vision-Audio-Language-Event Benchmark Towards Time-Aware Omni-Modal Perception of Long Videos

- ❑ **Summary:** This paper introduces LongVALE, a benchmark designed to evaluate comprehensive understanding in omni-modal models that integrate vision, audio, and language. Models pre-trained on LongVALE demonstrate strong performance in both omni-modal reasoning and fine-grained reasoning tasks.

- ❑ **Novelty:** LongVALE enables systematic evaluation of omni-modal perception and supports assessment across a wide range of reasoning abilities, including fine-grained and cross-modal understanding.

- ❑ **Thoughts:** LongVALE was constructed through the combination of existing models and automated processes. Despite recent progress, fine-grained understanding of visual data—such as images and videos—remains a significant challenge for current models.



(d) Audio-visual correlation types

## VIDHALLUC: Evaluating Temporal Hallucinations in Multimodal Large Language Models for Video Understanding

- ❏ **Summary:** This paper proposes VIDHALLUC, a dataset for evaluating temporal reasoning and temporal hallucinations in video large language models (Video LLMs). The dataset is automatically constructed through the process illustrated in the figure (left), with final verification performed by human annotators.
- ❏ **Novelty:** VIDHALLUC enables a comprehensive assessment of how current Video LLMs handle temporal hallucinations, revealing notable performance gaps compared to human reasoning.
- ❏ **Thoughts:** Temporal reasoning and fine-grained video understanding remain challenging areas, and VIDHALLUC highlights the limitations of existing models in capturing accurate temporal dynamics.

## VideoEspresso: A Large-Scale Chain-of-Thought Dataset for Fine-Grained Video Reasoning via Core Frame Selection

- ❏ **Summary:** This paper introduces VideoExpresso, a benchmark for evaluating fine-grained video understanding, automatically constructed using Chain-of-Thought reasoning combined with existing multi-modal language models. A model trained on VideoExpresso achieves state-of-the-art performance on several existing datasets.

- ❏ **Novelty:** Key contributions include a focus on fine-grained video recognition and the use of Chain-of-Thought reasoning for large-scale data generation—an approach that enhances both interpretability and training quality.

- ❏ **Thoughts:** Several benchmarks for fine-grained video recognition were presented at CVPR, reflecting growing interest in this area. Pretraining on fine-grained tasks appears essential to achieving high accuracy, and both fine-grained recognition and Chain-of-Thought reasoning remain open challenges.

LIMIT.LAB
https://limitlab.xyz/

## LION-FS: Fast & Slow Video-Language Thinker as Online Video Assistant

- ❏ **Summary:** This paper proposes LION-FS (Fast-Slow), a new method for online video assistants, particularly in the context of first-person view (FPV) videos. The Fast track provides rapid responses using high frame-rate processing, while the Slow track performs deeper, fine-grained reasoning by hierarchically processing more detailed features.

- ❏ **Novelty:** The core novelty lies in adapting the SlowFast architecture—originally used for action recognition—to the domain of online video assistants, enabling both speed and detailed understanding.

- ❏ **Thoughts:** While the SlowFast concept itself is not new, its application to online video assistance is practical and promising. The approach could be extended to other video recognition tasks, and explicitly incorporating detection into the Slow track could further enhance performance.

## Video-Panda: Parameter-efficient Alignment for Encoder-free Video-Language Models

- ❏ **Summary:** This paper introduces Video-Panda, an encoder-free video LLM architecture. Instead of using a conventional video encoder, the model employs a Spatio-Temporal Alignment Block (STAB) to extract fine-grained spatial and temporal features directly from video inputs. Video-Panda outperforms models like Video-LLaMA and Video-ChatGPT, while using only 45 million parameters for visual encoding.

- ❏ **Novelty:** The key innovation lies in removing the typically large video encoder module, resulting in a lightweight yet high-performing model that maintains strong accuracy with significantly fewer parameters.

- ❏ **Thoughts:** Video-Panda is notably compact but still performs competitively across multiple benchmarks. A more comprehensive evaluation showing how it achieves this performance would be insightful, especially compared to similarly sized models. It would also be valuable to explore how the model scales and how its performance evolves with larger configurations. The proposed STAB module appears somewhat complex and warrants further examination.

## StreamingT2V: Consistent, Dynamic, and Extendable Long Video Generation from Text

- ❏ **Summary:** This paper presents StreamingT2V, a text-to-video generation method structured in three stages, capable of generating videos around 2 minutes long. The model incorporates a Condition Attention Module to align current frames with recent ones and an Appearance Preserve Module to maintain global consistency over long sequences. Finally, long videos are generated in an autoregressive manner using randomized blending. The approach significantly outperforms prior methods such as OpenSora and OpenSoraPlan.
- ❏ **Novelty:** The improvements stem from engineering refinements and training strategies that substantially boost generation quality over existing baselines.
- ❏ **Thoughts:** The method has already gained notable attention. It raises interesting questions about how different camera settings (e.g., movement, FPV) may require distinct design considerations. It also remains to be seen whether a multi-stage approach is truly optimal. Training such a system likely involves significant complexity.

## EIDT-V: Exploiting Intersections in Diffusion Trajectories for Model-Agnostic, Zero-Shot, Training-Free Text-to-Video Generation

- ❏ **Summary:** This paper introduces EIDT-V, a training-free optimization method for improving existing image-based text-to-video generation models. By incorporating frame-wise descriptions and inter-frame motion prompts, EIDT-V adjusts generation in the latent space, enhancing both per-frame quality and temporal coherence. The method achieves strong performance without requiring model retraining.

- ❏ **Novelty:** Training-free approaches for text-to-video generation remain relatively underexplored. EIDT-V enables image-level control through textual prompts, offering a creative and flexible mechanism with potential across various applications.

- ❏ **Thoughts:** Using distinct prompts for each frame is a compelling idea and may also support temporal alignment. It would be interesting to explore whether this approach could be extended to video-native generation models beyond image-based frameworks.

## VideoDirector: Precise Video Editing via Text-to-Video Models

❏ **Summary:** This paper presents VideoDirector, a method for direct editing of text-to-video (T2V) generation models. The approach introduces separate guidance for spatial and temporal features, making it easier to edit each aspect independently. Additionally, it proposes multi-frame null-text optimization, enabling fine-grained temporal editing. The method achieves state-of-the-art performance.

❏ **Novelty:** The paper offers a thorough analysis of limitations in existing inversion-then-editing approaches to video editing, such as tightly coupled spatial-temporal features and complex layout entanglement. VideoDirector addresses these issues by explicitly decoupling spatial and temporal editing mechanisms.

❏ **Thoughts:** The challenge of spatial-temporal coupling in traditional T2V generation is well known. This work provides a clear and structured approach to mitigating it, offering a promising direction for more controllable and precise video editing.

## Uni4D: Unifying Visual Foundation Models for 4D Modeling from a Single Video

- ❏ **Summary:** This paper introduces Uni4D, a method for generating 4D scenes (3D geometry over time) directly from video. Uni4D leverages a visual foundation model to extract depth, motion tracking, and segmentation from video frames. Using these outputs, the system estimates camera poses and scene geometry, followed by bundle adjustment to refine results. Notably, the visual foundation model is used without retraining, yet Uni4D achieves top performance across multiple benchmarks.
- ❏ **Novelty:** The key innovation lies in applying a pre-trained visual foundation model to the video-to-4D task, enabling high-quality results without additional training or supervision.
- ❏ **Thoughts:** Though built from existing components, the method effectively combines them to deliver state-of-the-art performance. The integration of foundation models into dynamic scene reconstruction is promising and practical.

## LATTE-MV: Learning to Anticipate Table Tennis Hits from Monocular Videos

❏ **Summary:** This work makes two key contributions toward developing a table tennis agent. First, it proposes a pipeline to reconstruct the 3D scene from video, including mask extraction for the table and ball, estimation of the table's four corners, camera calibration, pose estimation, and final ball trajectory generation. Second, it introduces an uncertainty-aware controller based on generative modeling to predict the opponent's future actions. A dataset of approximately 50 hours of professional table tennis matches was also collected from online sources.

❏ **Novelty:** The inclusion of anticipation modeling—predicting opponent behavior under uncertainty—distinguishes this work from prior approaches.

❏ **Thoughts:** Incorporating physical dynamics would further enhance realism. Achieving full 4D scene reconstruction (geometry + motion over time) would be a promising next step.



**Algorithm 1** Ball Trajectory Reconstruction

**Require:** 2D ball positions $\{b_{2D,t}\}$, camera intrinsic parameters $\mathbf{K}$ and extrinsic parameters $\mathbf{R}, \mathbf{t}$
1: Find hit times $\{h_i\}_{i=1}^{H}$
2: **for** each $i = 1$ to $H - 1$ **do**
3:     Find potential bounce times $\{b_j\}_{j=1}^{B} \subset [h_i, h_{i+1}]$.
4:     **for** each $j = 1$ to $B$ **do**
5:         Fit parabolas to $\{b_{2D,t}\}_{t=h_i}^{b_j}$ and $\{b_{2D,t}\}_{t=b_j}^{h_{i+1}}$.
6:         Compute $\text{MSE}_j$ for each fit.
7:     **end for**
8:     Select $j^* \in \arg\min_{j \in [B]}\{\text{MSE}_j\}$ and set $b = b_{j^*}$.
9:     Set $b_{3D,h_i}$ to player's racket hand at frame $h_i$.
10:     Set $b_{3D,h_{i+1}}$ to player's racket hand at frame $h_{i+1}$.
11:     Compute $b_{3D,b}$ via inverse camera projection.
12:     Fit $b_{3D,t}$ for $t \in [h_i, h_{i+1}]$ via Eq. (2)–(5).
13: **end for**

Figure 7. Simulated target poses and ball trajectory.

## SnapGen-V: Generating a Five-Second Video within Five Seconds on a Mobile

- ❏ **Summary:** This paper proposes SnapGen-V, a lightweight yet high-performing text-to-video generation model capable of generating 5-second videos within 5 seconds on a mobile device (iPhone 16 Pro Max). To achieve this efficiency, the authors conducted extensive experiments to identify a compact yet effective temporal backbone, and designed efficient temporal layers. Additionally, a new adversarial fine-tuning strategy for diffusion models is introduced, along with distillation techniques that reduce the required denoising steps to just four.

- ❏ **Novelty:** The work stands out for its systematic exploration of lightweight architectures and presents the first text-to-video generation model that runs effectively on mobile devices.

- ❏ **Thoughts:** Model efficiency is an increasingly popular research direction. This work opens the door for real-time, on-device video generation, and may pave the way for a wave of mobile applications in the near future.

## Koala-36M: A Large-scale Video Dataset Improving Consistency between Fine-grained Conditions and Video Content

- ❏ **Summary:** This paper introduces Koala-36M, a large-scale dataset designed for both video recognition and generation. Each video clip (around 10 sec) is annotated with detailed textual descriptions averaging around 200 words, generated automatically by models. The authors also propose a Video Training Suitability Score (VTSS) to assess and enhance video quality. Furthermore, diffusion models are employed to improve temporal consistency between the video and its corresponding text.

- ❏ **Novelty:** Koala-36M provides more detailed and higher-quality captions compared to similar datasets such as Panda70M. The use of VTSS and diffusion models for quality and consistency further enhances its utility.

- ❏ **Thoughts:** High-quality, fine-grained datasets for video and image understanding are becoming increasingly essential. A key question is how to ensure the quality of automatically generated text. Still, generating large volumes of reasonably good captions with current models seems to be an effective strategy for training.



| Dataset | #Videos | ATL(words) | TVL(hours) | Text | Filtering | Resolution |
|---|---|---|---|---|---|---|
| LSMDC (Rohrbach et al., 2015) | 118K | 7.0 | 158 | Manual | Sub-metrics | 1080p |
| DiDeMo (Anne Hendricks et al., 2017) | 27K | 8.0 | 87 | Manual | Sub-metrics | - |
| YouCook2 (Zhou et al., 2018) | 14K | 8.8 | 176 | Manual | Sub-metrics | - |
| ActivityNet (Caba Heilbron et al., 2015) | 100K | 13.5 | 849 | Manual | Sub-metrics | - |
| MSR-VTT (Xu et al., 2016) | 10K | 9.3 | 40 | Manual | Sub-metrics | 240p |
| VATEX (Wang et al., 2019) | 41K | 15.2 | ~115 | Manual | Sub-metrics | - |
| WebVid-10M (Bain et al., 2021) | 10M | 12.0 | 52K | Alt-Text | Sub-metrics | 360p |
| HowTo100M (Miech et al., 2019) | 136M | 4.0 | 135K | ASR | Sub-metrics | 240p |
| HD-VILA-100M (Xue et al., 2022) | 103M | 17.6 | 760.3K | ASR | Sub-metrics | 720p |
| VidGen (Tan et al., 2024) | 1M | 89.3 | - | Generated | Sub-metrics | 720p |
| MiraData (Ju et al., 2024) | 330K | 318.0 | 16K | Generated & Struct | Sub-metrics | 720p |
| Panda-70M (Chen et al., 2024b) | 70M | 13.2 | 167K | Generated | Sub-metrics | 720p |
| **Koala-36M (Ours)** | 36M | 202.1 | 172K | Generated & Struct | Expert Model | 720p |

## AnomalyNCD: Towards Novel Anomaly Class Discovery in Industrial Scenarios

- ❏ **Summary:** This paper proposes AnomalyNCD, a new method for multiclass anomaly detection. Since anomaly regions typically occupy only a small portion of the image, the authors introduce a main element binarization (MEBin) module that emphasizes detecting central regions of anomalies. The method also incorporates mask-guided feature extraction to enrich semantic information through region segmentation. The approach achieves strong performance across several benchmarks.
- ❏ **Novelty:** The innovation mainly lies in the methodological design—while the individual components are not entirely novel, their integration leads to notably high accuracy.
- ❏ **Thoughts:** Applying multi-modal language models (MLLMs) to anomaly detection could be a promising direction, potentially enabling not just detection but also interpretation of anomaly types. There's also potential to unify anomaly detection and object counting within a single framework.

## Protecting Your Video Content: Disrupting Automated Video-based LLM Annotations

- ❏ **Summary:** This paper introduces two watermarking-style methods—Ramblings and Mute—to address the issue of video LLMs misinterpreting or over-accessing video content, thus aiming to protect data privacy. The Ramblings method steers Video LLMs toward incorrect responses by embedding misleading signals, while Mute focuses on reducing output quality by targeting the model's end-of-sequence (EOS) behavior. Both methods demonstrate strong performance against models such as Video-ChatGPT, Video-LLaMA, and Video-Vicuna.
- ❏ **Novelty:** The work addresses a new and important research problem: protecting data privacy in video content from video LLMs. The proposed methods introduce adversarial-style defenses specific to multimodal large language models.
- ❏ **Thoughts:** As video recognition capabilities of LLMs rapidly improve, the demand for privacy-preserving and adversarial techniques targeting MLLMs is likely to grow. This work highlights a critical emerging direction in responsible AI development.

LIMIT.LAB
https://limitlab.xyz/

## SMTPD:ANewBenchmarkfor Temporal Prediction of Social Media Popularity

❏ **Summary:** This paper proposes SMTPD, a new benchmark for evaluating social media video popularity, with a particular focus on the temporal dynamics of popularity. Compared to existing datasets, SMTPD captures how video popularity evolves over time.

❏ **Novelty:** The task itself is novel and socially relevant—social media popularity has significant implications for public influence and economic value, making it an important area for machine learning research.

❏ **Thoughts:** The topic of social popularity is compelling, with potential for exploring factors behind popularity, transferring popularity to other content, and more. However, since popularity is not static and can change due to real-world user interactions, it's unclear whether the task can be fully isolated as a learning problem. This work also reflects a broader trend of applying MLLMs to complex, real-world systems like social media.



(a) A social media post, also serving as a sample in SMTPD.

(b) Box plots of popularity scores over time.

| Dataset | Source | Category | Samples | Language | Prediction Type |
|---------|--------|----------|---------|----------|-----------------|
| Mazloom [35] | Instagram | fast food brand | 75K | English | single |
| Sanjo [42] | Cookpad | recipe | 150K | Japanese | single |
| TPIC17 [47] | Flickr | - | 680K | English | single |
| SMPD [48] | Flickr | 11 categories | 486K | English | single |
| AMPS [11] | YouTube | shorts | 13K | Korean | single |
| SMTPD (ours) | YouTube | 15 categories | 282K | over 90 languages | sequential |

## Video-3D LLM: Learning Position-Aware Video Representation for 3D Scene Understanding

- ❏ **Summary:** This paper proposes Video-3D LLM, a new multimodal large language model capable of 3D recognition using video. The method treats 3D understanding as a form of dynamic video analysis, aligning video representations with 3D representations. The model achieves strong performance across several 3D recognition benchmarks.
- ❏ **Novelty:** The key innovation lies in enabling 3D understanding from video, allowing an MLLM to perceive 3D structure without requiring explicit 3D input, by bridging video and 3D modalities.
- ❏ **Thoughts:** Recognizing 3D information from video alone is highly practical. This approach makes 3D perception more accessible in real-world scenarios, where dedicated 3D input is often unavailable. It enhances usability and opens up new application possibilities for MLLMs in spatial understanding.

## VideoGEM: Training-free Action Grounding in Videos

- ❏ **Summary:** This paper introduces VideoGEM, the first training-free video grounding method that localizes actions and objects in video without additional model training. It extends the image-based grounding method GEM to video. Based on the observation that different attention layers correspond to actions and objects, the method dynamically adjusts layer weights. It also generates prompts for actions, verbs, and objects from action labels, using these prompts to guide attention and extract corresponding attention maps.
- ❏ **Novelty:** This is the first approach to achieve video grounding without training, leveraging attention layer dynamics and prompt-based guidance to localize relevant content.
- ❏ **Thoughts:** The training-free nature of VideoGEM is highly appealing, especially for scalable or low-resource settings. Its flexibility with prompts could enable fine-grained, customizable grounding and support the automatic construction of video-based datasets, opening up broader applications.

## The Devil is in Temporal Token: High Quality Video Reasoning Segmentation

- ❏ **Summary:** This paper proposes VRS-HQ, a new method for video reasoning segmentation. Unlike previous methods that use a single [SEG] token to detect target objects across video frames, VRS-HQ emphasizes the importance of temporal localization. It introduces a temporal attention token ([TAK]) to identify which frames contain the target object. The method also integrates similarity-based fusion and frame selection mechanisms to enhance temporal reasoning.

- ❏ **Novelty:** The key innovation is the introduction of the [TAK] temporal token, enabling more precise temporal grounding of objects. Despite its simplicity, the method achieves strong performance across multiple datasets and tasks.

- ❏ **Thoughts:** The [TAK] token shows potential for broader use in other video reasoning tasks, and it raises interesting questions about whether similar temporal tokens could be applied to audio or other sequential modalities. The approach is practical and generalizable.

## MoManipVLA: Transferring Vision-language-action Models for General Mobile Manipulation

❏ **Summary:** This paper presents a framework for transferring large-scale Vision-Language-Action (VLA) models to mobile manipulation tasks in robotics. The proposed method predicts waypoint representations from VLA outputs, which generalize well across various manipulation tasks. It also incorporates motion planning for both arm and base actions, improving the accuracy of movement and manipulation in real-world settings. The system achieves mobile manipulation with minimal additional training.

❏ **Novelty:** The key innovation is the transfer of VLA models to real-world mobile manipulation, using waypoints as an effective and generalizable intermediate representation.

❏ **Thoughts:** There's a growing interest in applying MLLMs to robotics. This work raises important questions about what scene representations are most effective for real-world tasks—waypoints, scene graphs, or others. Integrating structured representations may further enhance precision and generalization in robotic planning.

## OpenING: A Comprehensive Benchmark for Judging Open-ended Interleaved Image-Text Generation (Oral)

❏ **Summary:** This paper introduces OpenING, a comprehensive benchmark designed to evaluate interleaved image-text generation in multimodal large language models (MLLMs). OpenING includes 5,400 human-annotated instances spanning 56 real-world tasks such as travel guidance, design, and brainstorming. To support reliable evaluation, the authors also propose IntJudge, a novel automatic judge model trained via a custom data pipeline, which achieves 82.42% agreement with human ratings, surpassing GPT-based evaluators. Experiments demonstrate that current interleaved generation models still fall short, and the paper provides insights to drive future advancements in the field.

❏ **Novelty:** Interleaved generation still remains less discussed.

❏ **Thoughts:** Dataset for detailed generation is still insufficient therefore important. Pretraining models with interleaved multimodal data could possibly improve model performance. In real-world user inputs, the data could be an incorporation of various data types, therefore it is important to carefully think about what kind of input data is important, less addressed, and fundamentally hard.



(a) Trends in Generative Model Development

(b) Benefits of Interleaved Image-Text Generation

LIMIT.LAB

140

https://limitlab.xyz/

## Molmo and PixMo: Open Weights and Open Data for State-of-the-Art Vision-Language Models (Oral)

- ❑ **Summary:** This paper introduces Molmo, a family of open-weight vision-language models trained entirely from scratch using a newly collected dataset suite called PixMo. Without relying on proprietary VLMs, Molmo achieves state-of-the-art performance among open models and even outperforms several leading proprietary systems such as Claude 3.5 Sonnet and Gemini 1.5 Pro. Molmo involves detailed caption, impressive grounding ability, detailed knowledge-based reasoning.
- ❑ **Novelty:** Openness of data and model. While Molmo is open-sourced, it obtained with state-of-the-art result even compared with closed source commercial models.
- ❑ **Thoughts:** Detailed recognition and grounding ability  is very import, with them several recent works including molmo obtained improved performance.

## Rethinking Vision-Language Model in Face Forensics: Multi-Modal Interpretable Forged Face Detector

- ❏ **Summary:** This paper proposes M2F2-Det, a multi-modal deepfake detection model that not only classifies images as real or fake but also provides detailed natural language explanations for its decisions, enhancing interpretability. The model introduces a bridge adapter to align the image encoder with the LLM-based explanation generator, and uses forgery prompt learning and layer-wise LF-Tokens to better capture subtle facial manipulations. Through multi-stage training, M2F2-Det achieves strong generalization to unseen forgeries and sets new benchmarks in both detection accuracy and explanatory quality.

- ❏ **Novelty:** The task of simultaneously generating score and interpretations (left image). Also, the model is complicated but new (right image).

- ❏ **Thoughts:** Faces, gestures, gazes are important in social interaction recognition, but they are not really well addressed in data side. Instead of score, the ability to ground regions to explain the reason for forged faces could be more interesting. The model is a bit complicated.

## Attention Distillation: A Unified Approach to Visual Characteristics Transfer

- ❏ **Summary:** The self-attention mechanism in diffusion models inherently captures image style and semantic structure. This paper proposes a unified framework that utilizes this attention information to transfer visual characteristics—such as style and texture—from a reference image to a newly generated one.

- ❏ **Novelty:** The technical novelty lies in the introduction of an *Attention Distillation Loss*, which is defined as the discrepancy between the attention map of an ideally stylized image and that of the generated image. This enables effective transfer of visual properties through image optimization in the latent space via backpropagation. Additionally, by integrating the attention distillation loss into the diffusion sampling process and enhancing the traditional classifier guidance, the method achieves faster and more efficient image generation.

- ❏ **Thoughts:** While the model excels in transferring style and texture, it still struggles with conveying higher-level semantic attributes such as composition and content. Nonetheless, due to its general applicability, the method holds strong potential for future development and practical utility—**especially in use cases that require consistent visual styling across images**.



https://github.com/xugao97/AttentionDistillation

## Estimating Body and Hand Motion in an Ego-sensed World

- ❏ **Summary:**This work proposes a system that estimates full-body 3D pose, body height, and hand motion in a global (world) coordinate system using only head-mounted device inputs—namely, 6DoF head poses from egocentric SLAM and onboard camera images. Crucially, the system **constructs the world map online**, without requiring any pre-scanned environments or external sensors.
- ❏ **Novelty:**The method introduces a **diffusion-based motion prior** conditioned on head movement. By performing **pose estimation in a locally normalized space relative to the head**, rather than in absolute coordinates, the model achieves significantly improved accuracy through spatial and temporal invariance.
- ❏ **Thoughts:**It is impressive that the method achieves full-body and hand motion estimation using only egocentric head pose and images. Height estimation is enabled under the assumption of a **flat ground plane**, which simplifies the problem. However, this also imposes a limitation: the method may not generalize well to **uneven terrains or stairs** due to this assumption.



There's a great video!

https://egoallo.github.io/

## Scene-Centric Unsupervised Panoptic Segmentation

❑ **Summary:** This paper proposes a novel framework for *scene-centric* ***unsupervised panoptic segmentation***, which enables holistic understanding of complex scenes *without any manual annotations*.

❑ **Novelty:** The method leverages stereo videos and self-supervised techniques (e.g., SMURF) to extract motion cues, followed by clustering via SF2SE3 to segment moving objects in 3D space. It further employs self-supervised features from DINO and Depth-G, incorporating stereo depth to generate high-resolution semantic masks. By combining motion-based instance masks ("things") with semantics-based region masks ("stuff"), the framework constructs high-quality pseudo panoptic labels.

❑ **Thoughts:** This method is highly promising for generating reliable pseudo labels even in complex real-world scenes. Its practicality is notable, as it can serve as an **effective pre-labeling step before manual annotation**, reducing annotation costs significantly.



https://visinf.github.io/cups/

## Structure-Aware Correspondence Learning for Relative Pose Estimation

- ❏ **Summary:**
- ❏ This paper addresses the task of estimating the 3D relative pose—rotation and translation—between a reference image and a query image, even when the object is previously unseen.
- ❏ **Novelty:**
- ❏ The method automatically extracts structure-representative keypoints from each image, allowing it to capture object shape structures even across significant appearance differences. By modeling the relationship between images as a graph, the approach directly predicts 3D-to-3D correspondences without requiring explicit feature matching.
- ❏ **Thoughts:**
- ❏ A key strength of this method lies in its ability to represent object structures as graphs and directly infer correspondences. However, since the structural estimation for a reference image changes depending on the query image, the method lacks a notion of absolute structure. On the other hand, this dependency allows it to reinterpret an image dynamically based on the query, making it practical for real-world applications such as object tracking or viewpoint adaptation.



https://cyhhzo02.github.io/SAC-Pose/

## Latent Drifting in Diffusion Models for Counterfactual Medical Image Synthesis

- ❏ **Summary:** Generating high-quality medical images is challenging due to the scarcity of labeled data. This paper introduces a novel method called **Latent Drifting (LD)** that enables counterfactual medical image synthesis by gradually shifting the latent distribution between source (e.g., natural) and target (e.g., medical) domains.
- ❏ **Novelty:** The key innovation lies in the introduction of a drift parameter δ into the denoising and reconstruction process of diffusion models. This allows controlled deviation in the latent space, facilitating domain adaptation. Furthermore, the method formalizes counterfactual generation as a min-max optimization problem—modeling "what-if" scenarios by modifying semantic attributes in a principled way.
- ❏ **Thoughts:** The approach is notable not only for bridging the distribution gap between natural and medical domains through latent control, but also for its structured treatment of counterfactual generation. This enables reliable medical image synthesis even in low-data regimes, with promising applications in clinical decision support and disease progression modeling.



https://latentdrifting.github.io/

## Multi-modal Vision Pre-training for Medical Image Analysis

- ❑ **Summary:** MRI's key strength is its ability to capture multiple modalities from the same anatomical cross-section, which inherently creates correlations between them. For effective self-supervised learning on MRI images, it's crucial to utilize datasets that include these diverse modalities.
- ❑ **Novelty:** To leverage these inter-modal correlations, BrainMVP constructed an impressive dataset: **16,022 scans (over 2.4 million images)** across **8 distinct MRI modalities**. They then conducted multi-modal image pre-training using three novel proxy tasks—**cross-modal image reconstruction, modality-aware contrastive learning, and modality template distillation**—to boost the learning of cross-modality representations and correlations, resulting in excellent performance on downstream tasks.
- ❑ **Significance:** The sheer scale of the dataset built for this project is truly remarkable. This dataset alone represents a significant contribution to medical imaging research.



https://github.com/shaohao011/BrainMVP

## Memories of Forgotten Concepts

- ❏ **Summary:** The paper reveals that erased concept information persists within AI models and can be revived. It shows that specific "latent seeds" and "inversion methods" can generate high-quality images of forgotten concepts.
- ❏ **Novelty:** This research uniquely demonstrates the intractability of completely erasing concept information from AI. It highlights potential vulnerabilities in current concept ablation techniques by showing how erased data can be reconstructed.
- ❏ **Thoughts:** I was surprised to learn that concepts exist in the latent space, not as discrete units, but as overlapping and dispersed information. This makes true "erasure" an incredibly challenging task, like trying to remove a specific color from a painting by only adjusting a few pixels.



https://matanr.github.io/Memories_of_Forgotten_Concepts/

## MotionPRO: Exploring the Role of Pressure in Human MoCap and Beyond

- ❏ **Summary:** This paper introduces MotionPRO, a large-scale human motion capture dataset integrating **pressure, RGB, and optical sensors**. It demonstrates the **necessity and effectiveness of pressure signals** for improving pose and trajectory estimation, and for driving realistic virtual humans.
- ❏ **Novelty:** The core novelty lies in creating a **multi-modal dataset** that rigorously explores the **critical role of pressure information** in human motion capture, thereby enhancing the physical realism and accuracy of AI models for motion understanding and generation.
- ❏ **Thoughts:** The inclusion of simultaneous pressure sensor data offers significant practical utility, especially for **time-series analysis of dynamic movements like jumps**. It provides crucial physical insights beyond visual cues, enabling more accurate and robust understanding of ground interaction and force distribution during such complex actions.



https://nju-cite-mocaphumanoid.github.io/MotionPRO/

LIMIT.LAB  150
https://limitlab.xyz/

## CRISP: Object Pose and Shape Estimation with Test-Time Adaptation

- ❏ **Summary:** CRISP is a new method that estimates 6D object pose and shape from a single RGB-D image. It excels at adapting to unseen objects and bridging large domain gaps through test-time self-training.
- ❏ **Novelty:** CRISP's primary novelty lies in its category-agnostic approach to pose and shape estimation, eliminating the need for prior object knowledge. It also introduces an effective test-time adaptation mechanism, allowing for strong generalization to novel objects and environments.
- ❏ **Thoughts:** It's genuinely impressive how CRISP manages to perform object pose and shape estimation without needing to know the object's category beforehand. This category independence is a significant leap forward, making it incredibly versatile for real-world scenarios with unknown items.

LIMIT.LAB

https://limitlab.xyz/

## Lessons and Insights from a Unifying Study of Parameter-Efficient Fine-Tuning (PEFT) in Visual Recognition

- ❏ **Summary:** This paper unifies the study of Parameter-Efficient Fine-Tuning (PEFT) methods in visual recognition, showing that various PEFT approaches achieve similar accuracy after careful tuning.
- ❏ **Novelty:** It highlights that despite similar performance, different PEFT methods make distinct predictions and errors, suggesting varying inductive biases.
- ❏ **Thoughts:** It's surprising to learn that each **PEFT method has a unique functional bias**, implying no single "one-size-fits-all" solution and necessitating method-specific applications.:



(a) Accuracy gain vs. linear probing on VTAB-1K (19 tasks) (b) Prediction overlaps (5K most confident) (c) Target distribution vs. distribution shifts

https://zheda-mai.github.io/PEFT_Vision_CVPR25/

## CraftsMan3D: High-fidelity Mesh Generation with 3D Native Diffusion and Interactive Geometry Refiner

❏ Summary: This work developed a new 3D generation system that can produce high-quality 3D shapes in a short time and allows interactive shape editing.

❏ Novelty:This work propose a generation model for 3D data that is fast, high-quality, and editable.

❏ Thoughts:I believe this is a state-of-the-art method capable of automatically generating 3D assets. If combined with LLM-based 3D scene generation, it could potentially solve the problem of the shortage of 3D data all at once.

## FirePlace: Geometric Refinements of LLM Common Sense Reasoning for 3D Object Placement

- ❏ Summary: This work investigate how to optimally leverage MLLMs for object placement tasks. This work proposed FirPlace.
- ❏ Novelty:This work show that presenting the options in multiple stages improves the inference accuracy of the MLLM.
- ❏ Thoughts:

## Thinking in Space: How Multimodal Large Language Models See, Remember, and Recall Spaces

- ❏ Summary: MLLMs show emerging spatial intelligence from video but struggle with reasoning, and explicit cognitive maps help
- ❏ Motivation & Objective: To investigate whether MLLMs can develop human-like spatial understanding from videos and identify key limitations in their visual-spatial reasoning
- ❏ Approach:
  - ❏ The study introduces VSI-Bench, a benchmark with diverse spatial tasks from real-world videos, and evaluates MLLMs by analyzing both their linguistic self-explanations and visual cognitive maps.
  - ❏ The dataset is built by unifying video and 3D annotation data from ScanNet, ScanNet++, and ARKitScenes, and generating over 5,000 QA pairs across eight spatial tasks using auto-annotation and human-in-the-loop refinement.

## Thinking in Space: How Multimodal Large Language Models See, Remember, and Recall Spaces

- ❏ Summary: Sonata mitigates the "geometric shortcut" problem in point cloud self-supervised learning by obscuring spatial cues and strengthening reliance on input features, leading to reliable and generalizable 3D representations that excel in linear probing and transfer tasks.
- ❏ Approach:
  - ❏ Macro Framework: Point Self-Distillation
  - ❏ Two-view strategy: Generate augmented point cloud views (crop, jitter, mask)
  - ❏ Align features: Use EMA teacher–student to match features between:
    - ❏ Local/masked views (student)
    - ❏ Global views (teacher)
  - ❏ Self-distillation replaces contrastive/generative learning
  - ❏ Loss: Sinkhorn clustering + KoLeo regularization for robust alignment

## CVPR AI arts

❏ AI art installations were exhibited in Hall A2. The most impressive piece for me was "The Flower."

❏ "The Flower" responds to human presence by tracking the viewer's face, and its petals move in a way that seems to convey emotion or intention — as if the flower were silently reacting to the observer.

## MoGe: Unlocking Accurate Monocular Geometry Estimation for Open-Domain Images with Optimal Training Supervision:

- ❏ Summary:This work propose an architecture that generates a 3D point cloud map from a single input image.
- ❏ Novelty: This work introduce an "affine-invariant" point map that is invariant to scale and translation.This removes focal-distance ambiguity and enables more stable learning,allowing direct 3D reconstruction from a single image without the need for camera parameters.
- ❏ Thoughts:I was particularly interested in how this could be applied to recognition tasks. Specifically, I wondered whether there might be a method to automatically generate annotations when reconstructing 3D scenes from a single image. For instance, if the input image already has semantic supervision or labels, it might be possible to propagate those annotations during the scene reconstruction process.



Input    Ground truth    Ours

## Multi-view Reconstruction via SfM-guided Monocular Depth Estimation

- ❏ Summary: This work propose a novel approach that incorporates prior information from SfM into diffusion-based depth estimation, enabling high-precision and multi-view consistent depth prediction for each viewpoint.
- ❏ Novelty: Using a diffusion model conditioned on SfM point clouds, we predict scale-consistent depth from multi-view images without the need for matching, and use this to achieve high-precision 3D reconstruction.
- ❏ Thoughts:

## Image Generation Diversity Issues and How to Tame Them

❏ Summary: This work proposes a new measurement method IRS for assessing diversity in image generation tasks. Using this metric, it is observed that conventional diffusion models—even when conditioned—are only able to cover up to 77% of the diversity present in the training data. Unconditional models perform even worse. Motivated by this insight, the authors propose a pseudo-labeling approach for generation, aiming to enhance diversity without sacrificing image quality.

❏ Novelty: This work defines new metrics of diversity and indicates conventional models remain low-diversity, even using conditional-generation.



Low Diversity / High Diversity



Predicted Diversity IRS$_\infty$ (%)

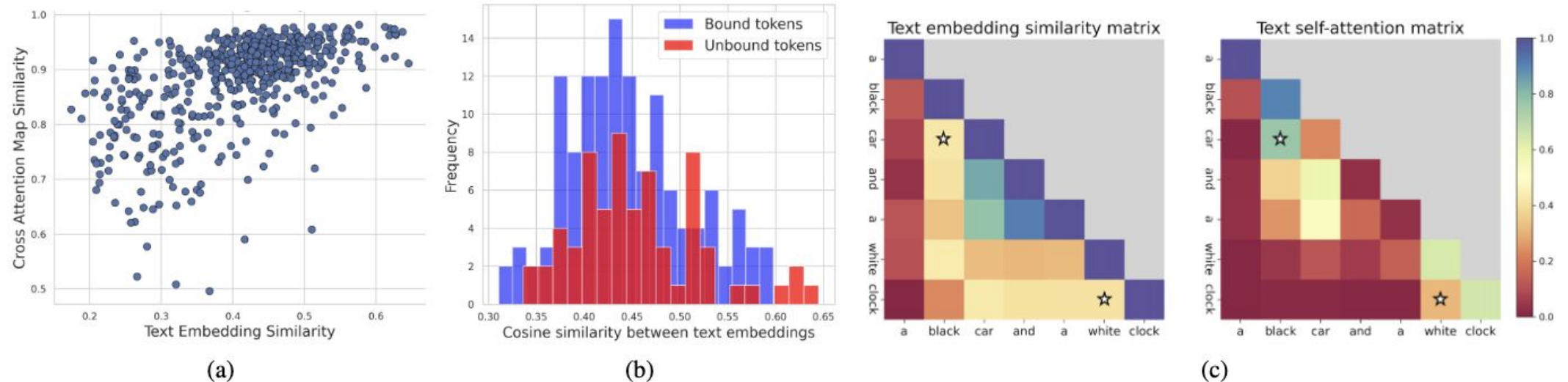| Model | Image Resolution | FID ↓ | Prec. ↑ | Rec. ↑ | Dens. ↑ | Cov. ↑ | Vendi ↑ | IRS$_{\infty,a}$ ↑ (Ours) |
|---|---|---|---|---|---|---|---|---|
| **Pixel diffusion** | | | | | | | | |
| ADM-256 [12] | 256 | 6.01 (30.30) | 0.82 (0.57) | 0.62 (0.73) | 1.08 (0.41) | 0.91 (0.40) | 70.94 (36.18) | 0.44 (0.20) |
| **Transformer** | | | | | | | | |
| DiT-XL/2-256 [36] | 256 | 22.15 (**8.72**) | 0.94 (0.69) | 0.34 (**0.76**) | 1.58 (0.70) | 0.85 (**0.84**) | 126.96 (**58.15**) | 0.23 (0.33) |
| DiT-XL/2-512 [36] | 512 | 22.99 (9.54) | **0.96** (0.70) | 0.27 (0.73) | **1.90** (0.72) | 0.86 (0.82) | **128.98** (55.17) | 0.21 (0.34) |
| MAR-B-256 [27] | 256 | 3.79 (10.36) | 0.83 (**0.72**) | 0.67 (0.71) | 1.18 (0.72) | 0.96 (0.75) | 83.03 (55.78) | 0.45 (**0.38**) |
| MAR-L-256 [27] | 256 | 3.30 (10.36) | 0.82 (**0.72**) | 0.71 (0.71) | 1.10 (0.73) | 0.96 (0.75) | 81.80 (55.95) | 0.56 (**0.38**) |
| MAR-H-256 [27] | 256 | 3.11 (10.36) | 0.82 (**0.72**) | **0.72** (0.71) | 1.07 (**0.74**) | 0.96 (0.76) | 81.37 (55.82) | 0.64 (**0.38**) |
| **Latent diffusion, U-Net** | | | | | | | | |
| LDM-256 [43] | 256 | 26.09 (37.39) | **0.96** (0.61) | 0.21 (0.68) | 1.80 (0.45) | 0.83 (0.28) | 126.94 (30.83) | 0.16 (0.16) |
| EDM-2-XS-512 [24] | 512 | 3.79 (75.02) | 0.83 (0.42) | 0.65 (0.63) | 1.22 (0.25) | 0.95 (0.13) | 72.41 (26.95) | 0.46 (0.09) |
| EDM-2-S-512 [24] | 512 | 3.33 (122.48) | 0.85 (0.33) | 0.67 (0.42) | 1.26 (0.16) | **0.97** (0.07) | 80.25 (21.34) | 0.59 (0.04) |
| EDM-2-M-512 [24] | 512 | 3.30 (107.45) | 0.85 (0.36) | 0.69 (0.61) | 1.24 (0.19) | **0.97** (0.09) | 82.99 (22.19) | 0.65 (0.06) |
| EDM-2-L-512 [24] | 512 | 2.90 (118.87) | 0.84 (0.23) | 0.70 (0.51) | 1.22 (0.11) | **0.97** (0.06) | 82.10 (22.89) | 0.71 (0.03) |
| EDM-2-XL-512 [24] | 512 | 2.92 (141.74) | 0.84 (0.25) | 0.71 (0.45) | 1.21 (0.12) | **0.97** (0.06) | 83.23 (20.04) | **0.77** (0.03) |
| EDM-2-XXL-512 [24] | 512 | **2.87** (124.29) | 0.84 (0.33) | 0.71 (0.60) | 1.22 (0.17) | **0.97** (0.07) | 82.45 (21.52) | 0.75 (0.05) |

## Text Embedding is Not All You Need: Attention Control for Text-to-Image Semantic Alignment with Text Self-Attention Maps

❏ Summary: This study investigated the issue where existing text-to-image models fail to faithfully follow the input text and proposed a solution. It identified that the problem arises because incorrect regions tend to be emphasized in the cross-attention map. The study revealed that, in text embeddings, tokens with high similarity are more likely to have high values in the cross-attention map.

❏ Novelty: Since the self-attention in the text encoder captures the structure of the text, the study improves faithfulness by encouraging the cross-attention to mimic the self-attention map through the exploration of latent variables.

## DepthCrafter: Generating Consistent Long Depth Sequences for Open-world Videos

❏ Summary: Depth estimation for open-world videos has been challenging due to the diversity in camera motion and content layout. Existing methods like Depth Anything-v2 lack temporal consistency when applied to videos. To address this, this paper proposes a method that leverages a pre-trained image-to-video diffusion model for temporally consistent depth estimation.

❏ Novelty: To handle diverse videos, the model is trained using real videos for content diversity and synthetic videos for accurate depth supervision. Additionally, for precise conditioning, it leverages not only cross-attention but also conditioning in the VAE latent space. Moreover, the model adopts a training strategy that selectively fine-tunes either temporal or spatial layers, enabling it to handle long-frame videos efficiently.



Figure 2. Overview of our *DepthCrafter*. It is a conditional diffusion model that models the distribution $p(\mathbf{d}\mid\mathbf{v})$ over the depth sequence $\mathbf{d}$ conditioned on the input video $\mathbf{v}$. We train the model in three stages, where the spatial or temporal layers of the diffusion model are progressively learned on our compiled realistic or synthetic datasets with variable lengths $T$. During inference, given an open-world video, it can generate temporally consistent long depth sequences with fine-grained details for the entire video from initialized Gaussian noise, without requiring any supplementary information, such as camera poses or optical flow.



Figure 1. We innovate DepthCrafter, a novel video depth estimation approach, that can generate temporally consistent long depth sequences with fine-grained details for open-world videos, without requiring additional information such as camera poses or optical flow.

LIMIT.LAB

https://limitlab.xyz/

## Trends of Gaussian Splatting

❏ 93 papers that contain "Gaussian Splatting" in the titles in 2025 vs 28 in 2024 (≈3.3× growth).

    ❏ Gaussian splatting is attracting attention as 3D scene reconstruction is one of the areas of focus at CVPR 2025.

❏ Comparison with CVPR 2024

    ❏ Efficiency / sparse methods: 26 papers (fastest-growing)

    ❏ Scene reconstruction & dynamic scenes ≈ tripled

    ❏ Human-avatar share shrank from 32% → 12%

❏ Gaussian Splatting shifts from niche to default for real-time 3D

    ❏ Replacing NeRF as the dominant 3D scene reconstruction method.

## NexusGS: Sparse View Synthesis with Epipolar Depth Priors in 3D Gaussian Splatting

- ❏ Summary: This work embeds epipolar depth priors into 3D Gaussian Splatting to improve sparse-view novel view synthesis.
- ❏ Novelty: It introduces a point cloud densification strategy initializing dense Gaussians with epipolar depth priors via optical flow, employing Flow-Resilient Depth Blending and Flow-Filtered Depth Pruning to suppress flow errors and produce accurate depth under sparse views.
- ❏ Thought: NexusGS cleverly uses epipolar geometry to overcome the limitations of sparse views, providing robust depth initialisation without the need for heavy regularisation. Its effective three-step process relies on optical flow, which may struggle in textureless regions.



[Y. Zheng et al, CVPR 2025.][Link]

## Dr. Splat: Directly Referring 3D Gaussian Splatting via Direct Language Embedding Registration

- ❏ Summary: Dr. Splat directly registers CLIP embeddings to 3D Gaussians, bypassing rendering for efficient open-vocabulary 3D scene understanding.
- ❏ Novelty: It introduces a direct feature registration technique that assigns CLIP embeddings to dominant 3D Gaussians and integrates pretrained Product Quantization for compact embedding representation without per-scene optimization.
- ❏ Thought: I find the idea of eliminating the rendering stage compelling. It reduces feature distortion and speeds up 3D queries. Pretrained PQ balances memory and accuracy effectively.



[K. Jun-Seong et al., CVPR 2025.][Link]

## Creating Your Editable 3D Photorealistic Avatar with Tetrahedron-constrained Gaussian Splatting

❏ Summary: A hybrid TetGS-based pipeline generates editable photorealistic 3D avatars from monocular videos with text/image-guided edits.

❏ Novelty: This work proposes TetGS: embedding Gaussian kernels in tetrahedral grids for decoupled spatial adaptation and appearance learning, enabling precise localized geometry edits and photorealistic texture generation with coarse-to-fine and few-shot supervision.

❏ Thought: This paper addresses the issue of instability in 3D Gaussian Splatting editing by integrating structured tetrahedral grids. This approach offers clear geometric control and high-fidelity results. Although computationally intensive, the paper's staged pipeline and hybrid representation represent a significant advance in the creation of practical, user-friendly 3D avatars.



(d) Image-guided editing

[H. Liu et al., CVPR 2025.][Link]

## From Photorealism to Geometric Integrity

- ❏ 3DGS weakness: geometrically inaccurate surfaces/blurry edges/multi-view inconsistency.
    - ❏ The primary research trend has decisively shifted from simply improving visual quality to correcting these geometric flaws
- ❏ The "Splat" is No Longer Just 3D blob (Gaussian).
    - ❏ To achieve geometric accuracy, a wave of new primitives has emerged.
        - ❏ **2D Surfel Splatting** for surface consistency
        - ❏ **3D Convex Splatting** for sharp-edged objects
        - ❏ **Deformable Beta Splatting** for higher fidelity with fewer parameters
- ❏ Integration with Graphics Pipelines
    - ❏ While advanced primitives were introduced, they were still custom representations.The ultimate goal for many applications, such as games and VFX, remains generating standard triangular meshes
        - ❏ **Triangle Splatting** represents a potential convergence point for the field, making the fundamental graphics primitive directly optimizable.



Tangent frame (u,v)

New primitives

2D Gaussian Splat in object space

Corner

1 Smooth Convex

For Applications

## RipVIS: Rip Currents Video Instance Segmentation Benchmark for Beach Monitoring and Safety

- ❏ **Summary** : Creating large-scale video instance segmentation benchmark (RipVIS) for detecting rip currents. (Author said it took around 3 years to complete this)

- ❏ **Comment** : Author said Japan already has a nice system to predict rip current but it is closed. So they want to find Japanese collaborator. A competition based on RipVIS will be held in ICCV 2025.



- ❏ Project page
- ❏ Rip Current Instance Segmentation Challenge  (ICCV 2025)

**Workshop:** [Domain Generalization: Evolution, Breakthroughs and Future Horizon](#)

❏ Speaker: Aditi Raghunathan

Title: Predicting the Performance of Foundation Models Under Distribution Shift

  ❏ OOD accuracy is almost linearly correlated with ID accuracy.

  ❏ ID and OOD agreement also correlate linearly iff ID and ODD accuracy do.

    ❏ This empirical phenomena will make possible to **estimate FM performance under distribution shift with unlabeled data** . ([paper](#))

  ❏ But computing agreement over multiple FMs is not computationally feasible.

    ❏ **Random initializations of linear heads** works (no need for multiple checkpoint).

Workshop: [Visual Generative Modeling: What's After Diffusion?](...)

❏ Speaker: Kaiming He

Title: Towards End-to-End Generative Modeling

❏ Before AlexNet layer-wise training was a popular solution in recognition.

❏ Today's generative models are still conceptually like "layer-wise training".

❏ Therefore, there may be significant room for improvement **if the generative model can be made end-to-end** .

❏ Approaches based on **Flow** appear to be a promising direction at present.

❏ Latest work: [Mean Flows for One-step Generative Modeling](...)

## Unseen Visual Anomaly Generation

❏ Summary:

  ❏ Proposed model (AnomalyAny) generates realistic and diverse anomaly images at test time using only normal samples and a pre-trained Stable Diffusion model.

❏ Novelty:

  ❏ It uniquely leverages cross-attention-guided optimization and prompt refinement to steer anomaly generation without any additional training.

❏ Thoughts:

  ❏ The zero-training approach is elegant, though incorporating real one-shot anomaly exemplars could further boost generation precision.



Figure 2. **Illustration of AnomalyAny** with details of the attention-guided & prompt-guided optimization process at time step $t$.

## Workshop on Video Large Language Models

- ❏ Evaluation benchmark
  - ❏ Hallucination and Omission (ARGUS)
  - ❏ Compositional reasoning (VELOCITI)
  - ❏ Road event (RoadSocial)
  - ❏ Temporal understanding (Lost in Time)
- ❏ Efficient handling of long videos
  - ❏ Multi-frame fusion (FiLA-Video)
  - ❏ State-space model (BIMBA)
  - ❏ Frame sampling (Moment Sampling)
- ❏ Video grounding
  - ❏ Weakly Supervised (CoSPaL, STPro)
  - ❏ Low cost (NumPro)
  - ❏ Pixel grounding (PG-Video-LLaVA)
- ❏ Others
  - ❏ Current bottlenecks in training
  - ❏ Zero-shot Action Localization

## BIMBA: Selective-Scan Compression for Long-Range Video Question Answering

❏ Summary

  ❏ Video understanding with MLLMs face considerable difficulties when dealing with long videos, primarily due to the prohibitive quadratic computational cost of the self-attention mechanism, which is central to processing vast numbers of spatiotemporal tokens.

❏ Novelty

  ❏ BIMBA uses Mamba's efficient selective scan to extract key information from high-dimensional video, reducing it to a compact, information-rich token sequence.

❏ Thoughts

  ❏ The "Vanilla" (linear, Q-former) method quickly runs into GPU memory issues, and pooling struggles to capture long-range dependencies. Therefore, SSM-based methods, which can efficiently model temporal information, are expected to advance further for video understanding.



The woman's primary goal is to buy groceries. she is shopping for items such as cereal, milk, and cookies.

Large Language Model

Selective-Scan Spatiotemporal Token Selector

(Optional) Text Encoder

Image Encoder | Image Encoder | Image Encoder | Image Encoder

Based on the actions performed by the woman in the video, determine her primary goal and support your conclusion with relevant evidence from the video.

Frame 1 | Frame 2 | Frame 3 | Frame N | Language Query

Output Projection / Selective Scan / Input Projection / Layer Norm

(a) Spatiotemporal Token Selector

(b) Selective Scan with queries appended.

(c) Selective Scan with interleaved queries.

(d) Bidirectional Selective Scan with interleaved queries.

Spatiotemporal Tokens | Vision Queries

Forward Scan | Backward Scan

## Watermark

- ❏ For generative model
  - ❏ [Robust Watermark against Fine-Tuning](#)
  - ❏ [Box-Free Watermark Removal](#)
  - ❏ [Black-Box Forgery Attacks on Semantic Watermarks for Diffusion Models](#)
- ❏ Localized watermark/Partial theft
  - ❏ [Robust Watermarking Scheme Against Partial Image Theft](#)
  - ❏ [Manipulation Localization](#)
- ❏ For 3D Gaussian Splatting
  - ❏ [GuardSplat](#), [3D-GSW](#)
- ❏ Others
  - ❏ [Robust Message Steganography](#)
  - ❏ [Open-source Dataset Copyright](#)

## Video LLM

- ❏ Frame/Token Selection and Compression
  - ❏ [Dynamic Compression of Tokens](#)
  - ❏ [Adaptive Keyframe Sampling](#), [Flexible Frame Selection](#), [M-LLM Based Video Frame Selection](#)
- ❏ Temporal understanding
  - ❏ [Sequential Knowledge Transfer](#)
  - ❏ [Consistency of Video Large Language Models in Temporal Comprehension](#)
  - ❏ [Fine-Grained Compositional and Temporal Alignment](#)
- ❏ Spatial-Temporal Grounding
  - ❏ [VideoRefer Suite](#), [LLaVA-ST](#), [VideoGLaMM](#)
- ❏ Others
  - ❏ [Mitigating Action-Scene Hallucination](#)
  - ❏ [Real-time interactive streaming](#)

## Hand-Object Interaction

- ❏ Dataset
  - ❏ HD-EPIC: A Highly-Detailed Egocentric Video Dataset
  - ❏ EgoPressure: A Dataset for Hand Pressure and Pose Estimation in Egocentric Vision
- ❏ 3D Generation
  - ❏ EasyHOI: Unleashing the Power of Large Models for Reconstructing Hand-Object Interactions in the Wild
  - ❏ LatentHOI: On the Generalizable Hand Object Motion Generation with Latent Hand Diffusion
  - ❏ HOIGPT: Learning Long-Sequence Hand-Object Interaction with Language Models

## Workshop :  4D vision Modeling the Dynamic World

1.  4D Gaussian Splatting
    4D LangSplat: 4D Language Gaussian Splatting via Multimodal Large Language Models 4D LangSplat
    TRL-GS: Cascaded Temporal Residue Learning for 4D GS
    TRL-GS

2.  3D Rigging, 4D object intrinsics generation
    CVPR Poster Category-Agnostic Neural Object Rigging
    Birth and Death of a Rose
    Anymate: A Dataset and Baselines for Learning 3D Object Rigging

3.  4D Reconstruction
    St4RTrack: Simultaneous 4D Reconstruction and Tracking in the  World  St4rtrack

LIMIT.LAB
https://limitlab.xyz/

## Insight : 4D vision Modeling the Dynamic World

・Gaussian Splatting is also powerful in 4D 4D LangSplat

- 4D LangSplat embeds multimodal-LLM-generated, temporally consistent object captions into a dynamic 4D Gaussian-Splatting scene representation, enabling open-vocabulary and time-aware semantic querying.

-

・Temporal object intrinsics generation

- The rose of Death learn temporal object intrinsics by distilled it's 4D prior from a 2d diffusion.

Birth and Death of a Rose

## VideoGigaGAN: Towards Detail-rich Video Super-Resolution

❏ **Summery**

    ❏ The proposed model (Video-GAN) achieves both high quality and temporal consistency across output frames, delivering higher detail than existing models.

❏ **Novelty**

    ❏ VideoGigaGAN can produce video results with both high-frequency details and temporal consistency while artifacts like aliasing are significantly mitigated.

❏ **Thoughts**

    ❏ In video super-resolution, the issue of achieving both frame-to-frame consistency and quality is likely to continue to be debated.



Input      BasicVSR++      GigaGAN      Ours

Yiran Xu et al. "VideoGigaGAN: Towards Detail-rich Video Super-Resolution", in CVPR 2025.

LIMIT.LAB

## SAMWISE: Infusing Wisdom in SAM2 for Text-Driven Video Segmentation(Highlight)

❏ **Summary:**
  ❏ SAMWISE leverages SAM2 for text-driven video segmentation by integrating natural language understanding and temporal modeling without fine-tuning, while also introducing a novel module to adjust SAM2's inherent tracking bias for improved performance on dynamic objects.
❏ **Novelty:**
  ❏ Efficiently achieves text-driven video segmentation without re-training the foundation model, while overcoming SAM2's tracking limitations.
❏ **Thoughts:**
  ❏ Similar to last year, we're still seeing a large number of derivative studies building upon SAM and SAM2. I think the integration of text prompts without fine-tuning is particularly impressive.

# LIMIT.Lab

A collaboration hub

for building multimodal AI models under limited resources

## Our Members

**Hirokatsu Kataoka**
AIST / Oxford VGG

**Yoshihiro Fukuhara**
AIST

**Rintaro Yanagi**
AIST

**Ryousuke Yamada**
AIST

**Daichi Otsuka**
AIST

**Partha Das**
UvA

**Nakamasa Inoue**
Science Tokyo

**Go Irie**
TUS

**Rio Yokota**
Science Tokyo

**Ikuro Sato**
Science Tokyo

**Rei Kawakami**
Science Tokyo

**Christian Rupprecht**
Oxford VGG

**Iro Laina**
Oxford VGG

**Yuki M. Asano**
UTN FunAI Lab

**Elliott (Shangzhe) Wu**
Cambridge

**Daniel Schofield**
Oxford VGG

# Limited Resources, Unlimited Impact with Multimodal AI Models

AI foundation models are increasingly dominating various academic and industrial fields, yet the R&D of related technologies is controlled by limited institutions capable of managing extensive computational and data resources. To counter this dominance, there is a critical need for technologies that can develop practical AI foundation models using the standard computational and data resources. It is said that the scaling laws no longer provide the reliable roadmap for developing AI foundational models. Our community (LIMIT.Community) and the international lab (LIMIT.Lab) therefore aim to put in place exactly those technologies that permit the construction of {Vision, Vision-Language, Multimodal}AI foundational models even when compute and data are limited. Drawing on our members' prior successes in (i) generative pre-training methods that can be applied horizontally across any modality with image, video, 3D, & audio, and (ii) high-quality AI models from extremely scarce data (including a single image), we have been committed to AI multimodal foundational models under very limited resources. As of 2025, LIMIT.Lab is composed primarily of international research teams from Japan, UK, and Germany. Through collaborative research projects and the workshop organization, we actively foster global exchange in the field of AI and related areas.

Left: Core members / Right: Our mission

LIMIT.LAB
https://limitlab.xyz/

# Representation Learning with Very Limited Resources: When Data, Modalities, Labels, and Computing Resources are Scarce

ICCV 2025 Workshop

October, 2025      Honolulu, Hawaii
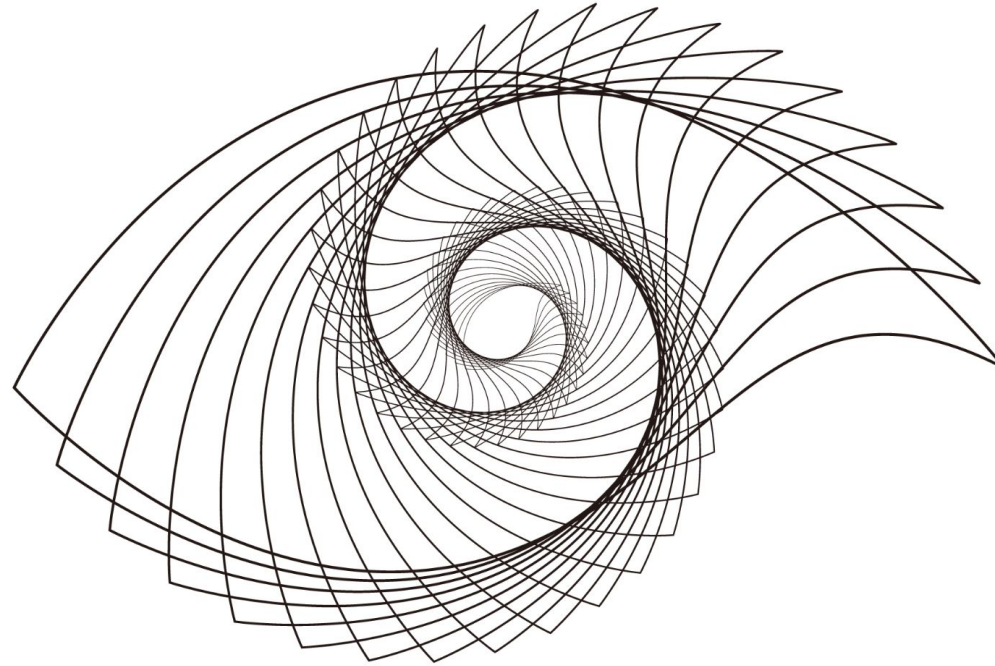
Submit Paper      Check Program

Please submit your paper to the ICCV25 LIMIT Workshop!

Deadline: July 10 (HST) / Length: 4 pages

## About LIMIT Workshop

Modern vision and multimodal models depend on massive datasets and heavy compute, magnifying costs, energy use, bias, copyright, and privacy risks. The "DeepSeek shock" of January 2025 spotlighted the urgency of learning powerful representations under tight resource limits. Now in its third edition, our workshop continues to explore strategies for robust representation learning

183

LIMIT.LAB

https://limitlab.xyz/

Join us! -> Slack invitation [Link]