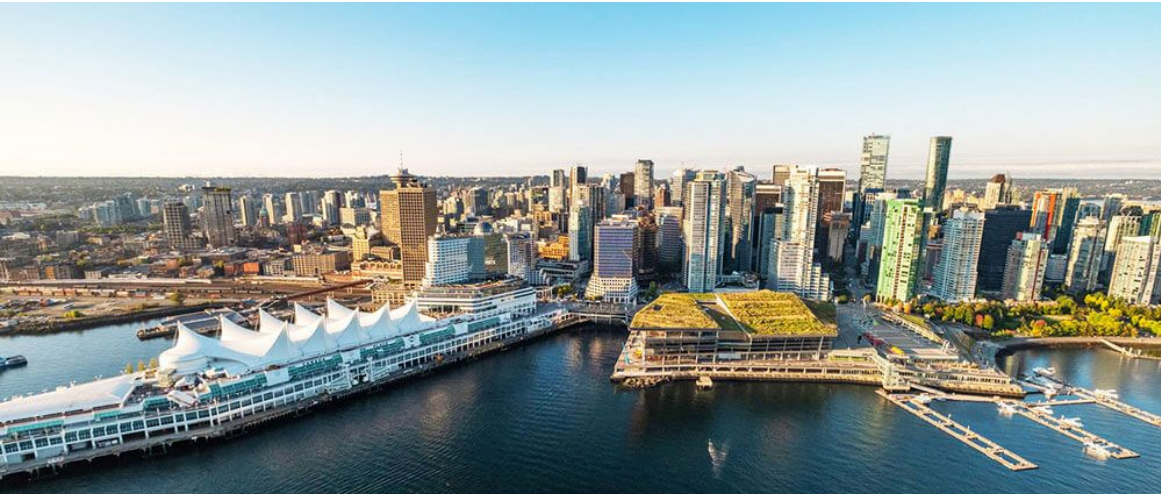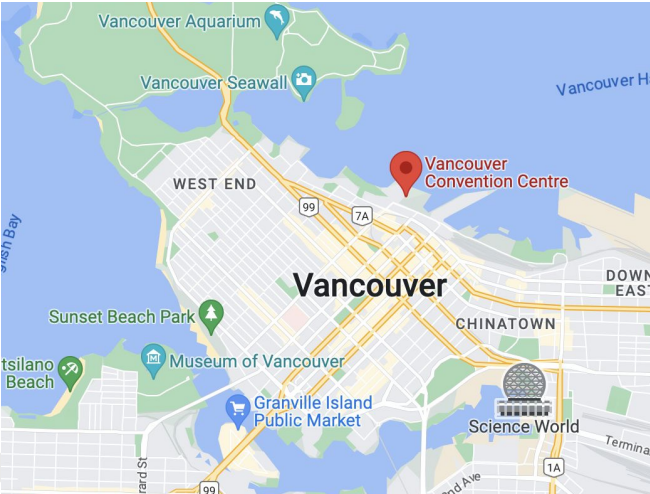**CVPR Impressions**

Papers

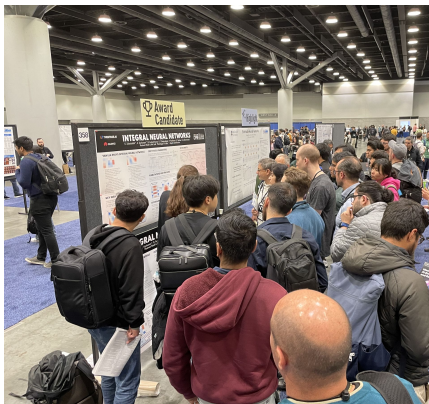Food for Thought

# Vancouver

# CVPR 2023

- Papers Submitted 9155
- Papers Accepted 2360
- Acceptance Rate 25.78%

- Attendance CVPR 23: 7088 in-person, 3215 virtual
- Attendance CVPR 19: 9375 (pre-COVID)

- Companies at CVPR 23: 116, 21200 square feet expo
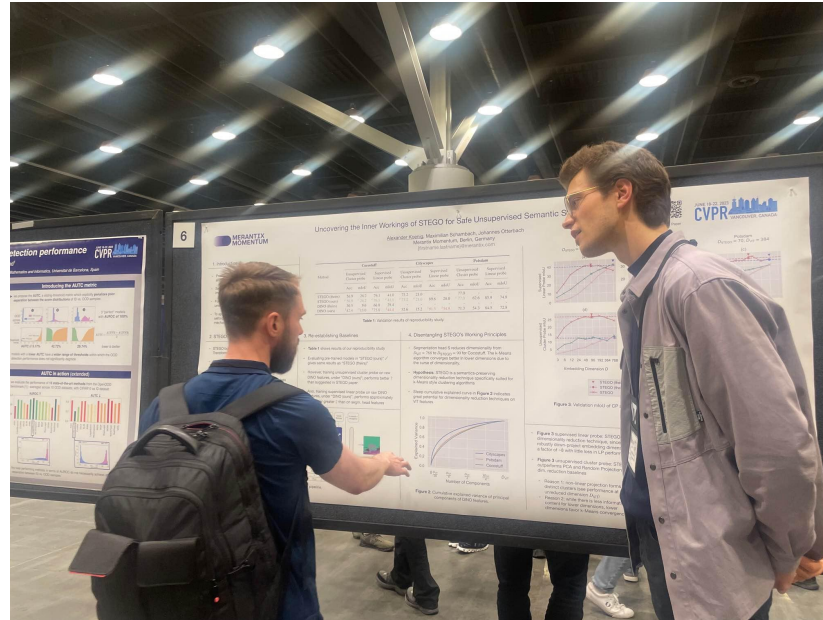- Companies at CVPR 19: 181, 41200 square feet expo
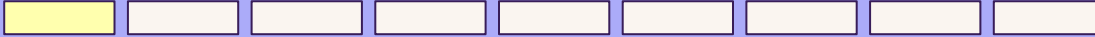
CVPR Impressions

**Papers**

Food for Thought

# Disclaimers

- Not my work! Just want to share some "aha" moments
- Hope to convey the paper's message
- If you want to deep dive, read the paper

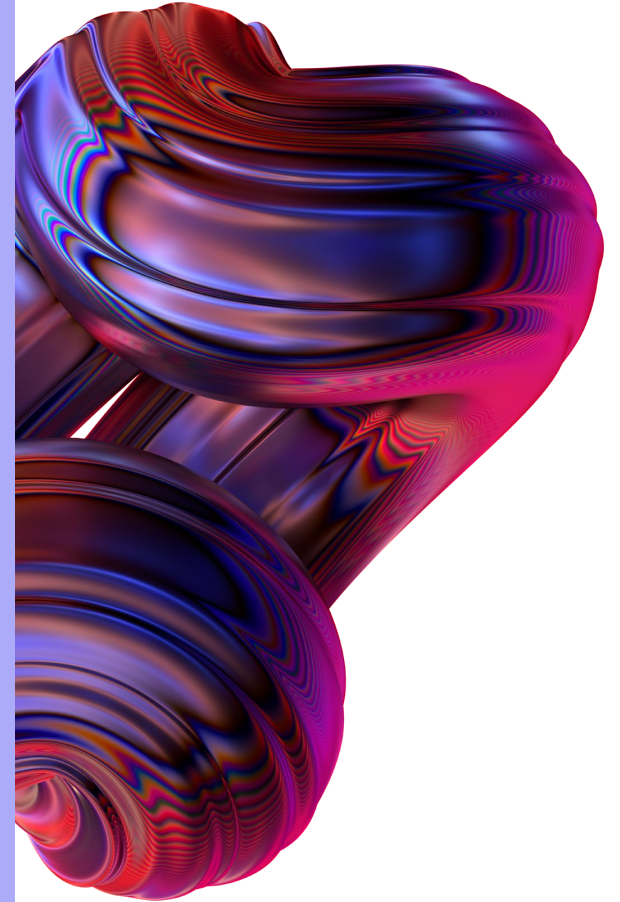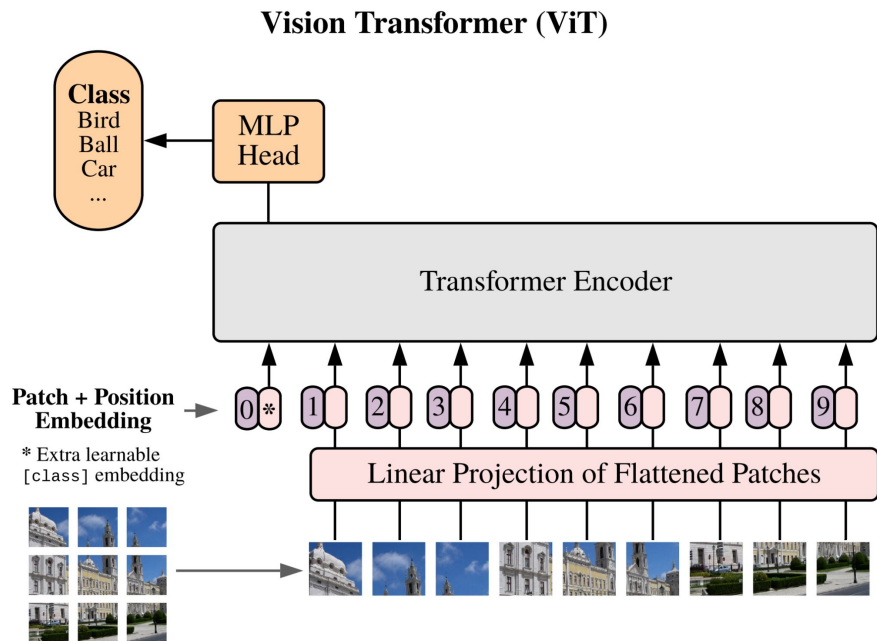# FlexiViT: One Model for All Patch Sizes

Lucas Beyer[*][1]    Pavel Izmailov[*][1,3]    Alexander Kolesnikov[*][1]    Mathilde Caron[*][2]    Simon Kornblith[*][1]

Xiaohua Zhai[*][1]    Matthias Minderer[*][1]    Michael Tschannen[*][1]    Ibrahim Alabdulmohsin[*][1]    Filip Pavetic[*][1]

Google Research

# Vision Transformers 101

**Vision Transformer (ViT)**



| | ViT-B/16 | ViT-B/32 | ViT-L/16 | ViT-L/32 |
|---|---|---|---|---|
| CIFAR-10 | 98.13 | 97.77 | 97.86 | 97.94 |
| CIFAR-100 | 87.13 | 86.31 | 86.35 | 87.07 |
| ImageNet | 77.91 | 73.38 | 76.53 | 71.16 |

Family of ViT Models

- **Problem**: need to train one model for each patch size (expensive, inflexible, must scale image s.t. 16 or 32 are a factor of resolution)

- **Trade-off**: small patch size → high performance, but expensive compute, and vice versa for large patch sizes

Source: Dosovitskiy et al. ICLR 2021

8

# FlexiViT - Key Idea



*Few Tokens*

Patch Embedding Weights

*Many Tokens*

FlexiViT — *Shared Weights* — FlexiViT

*84.4% INet Acc.*
*1.6 ms/img*

*86.1% INet Acc.*
*13 ms/img*

**Algorithm 1** Minimal FlexiViT pseudo-implementation.

```
1  model = ViT(...)
2  for batch in data:
3      ps = np.random.choice([8, 10, ..., 40, 48])
4      logits = model(batch["images"], (ps, ps))
5      # [...] backprop and optimize as usual
6
7  class ViT(nn.Module):
8      def __call__(self, image, patchhw):
9          # Patchify, flexibly:
10         w = self.param("w_emb", (32, 32, 3, d))
11         b = self.param("b_emb", d)
12         w = resize(w, (*patchhw, 3, d))
13         x = conv(image, w, strides=patchhw) + b
14         # Add flexible position embeddings:
15         pe = self.param("posemb", (7, 7, d))
16         pe = resize(pe, (*x.shape[1:3], d))
17         return TransformerEncoder(...)(x + pe)
```

**Notes**: Changes to existing code highlighted via violet background.

→ bilinear interpolation to resize patch embedding weights and positional embeddings
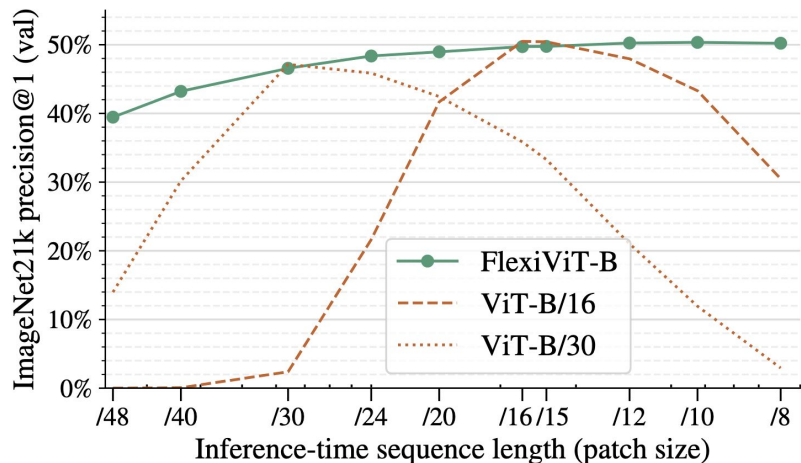
# FlexiViT - Results



Figure 3. **Standard ViTs are not flexible** in patch size. However, FlexiViT can train them to be flexible without loss of performance.
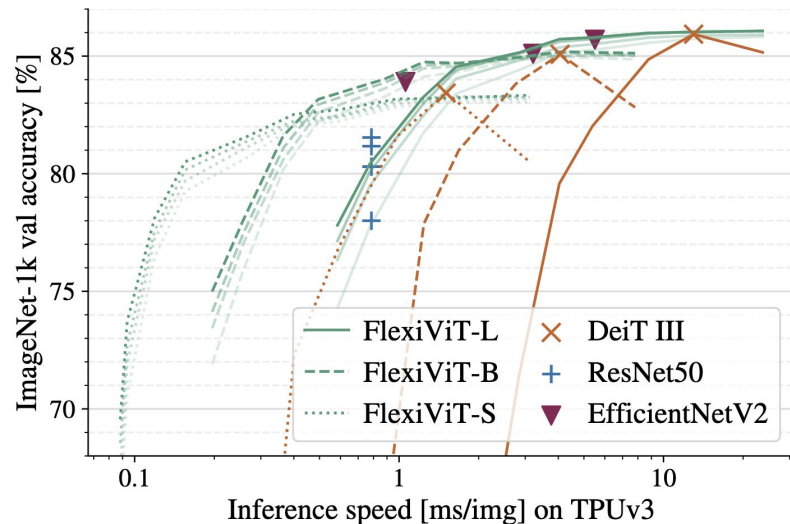


Figure 2. **FlexiViT results on ImageNet-1k.** We train three Flexi-ViTs based on DeiT III on ImageNet-1k and show their speed-accuracy tradeoff when evaluated at various patch sizes.

**Heuristic**: choose smallest patch size that still fits your compute budget ;-)

# Uncovering the Inner Workings of STEGO for Safe Unsupervised Semantic Segmentation
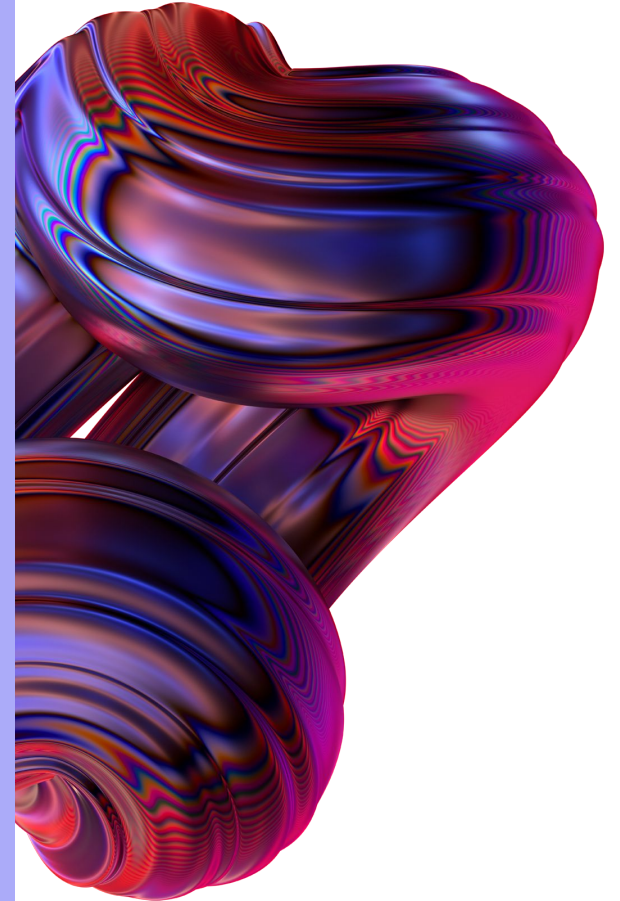
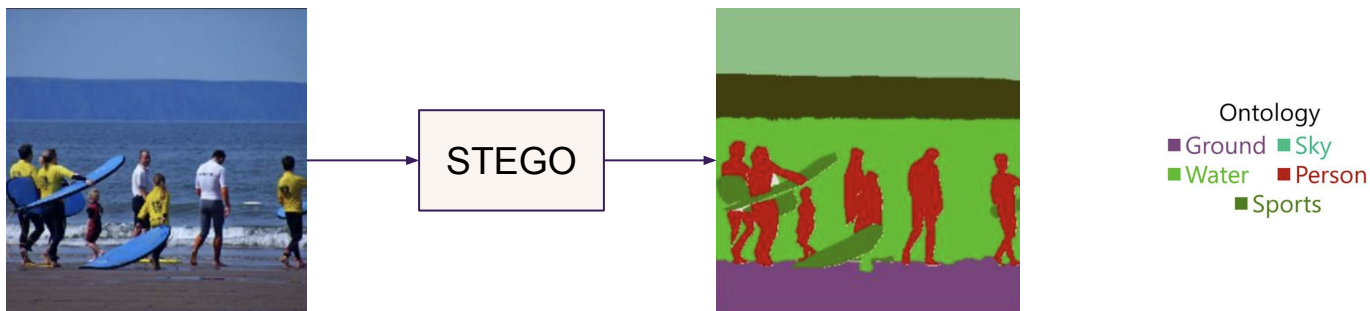Alexander Koenig          Maximilian Schambach          Johannes Otterbach

Merantix Momentum
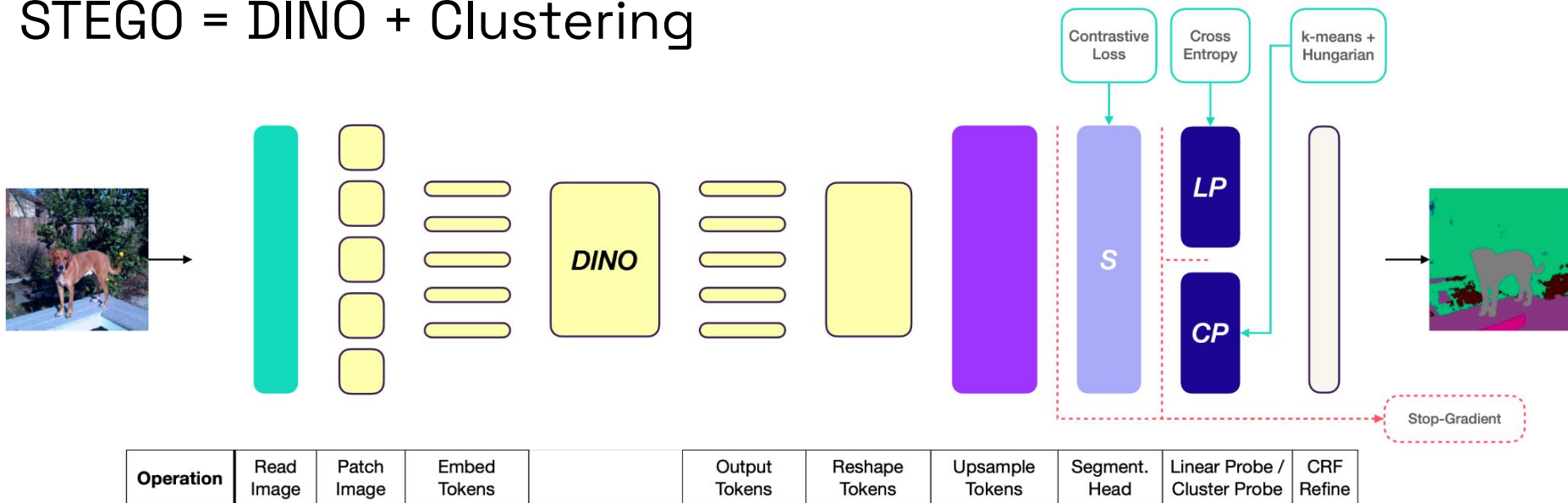
{firstname.lastname}@merantix.com

# STEGO Follow-Up: Motivation

- Problem: labeled data is scarce, but unlabeled data is abundant
- Self-supervised learning recently demonstrated impressive results on unlabeled datasets
- STEGO (Hamilton et al., ICLR 2022) does unsupervised semantic segmentation
- To apply STEGO safely in real-world, it's crucial to understand its working mechanisms
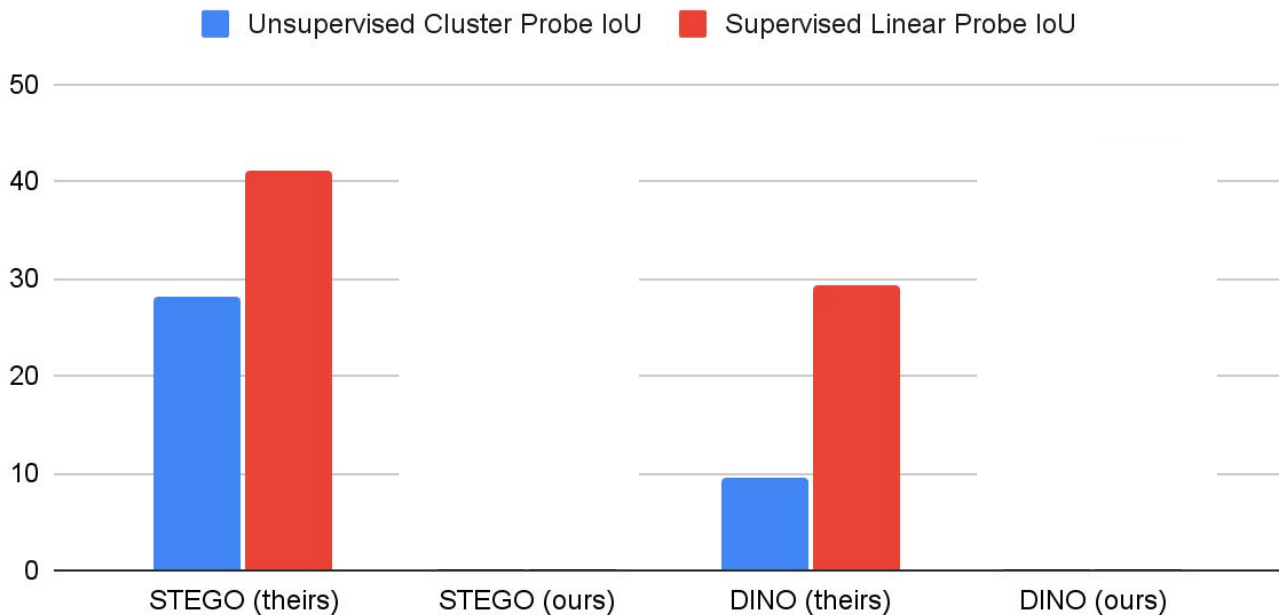
# STEGO = DINO + Clustering



- STEGO builds on DINO (Caron et al., ICCV 2021) pre-trained Vision Transformer
- Segmentation head S projects DINO feats into lower-dimensional space, "distilling" DINO feature correspondences
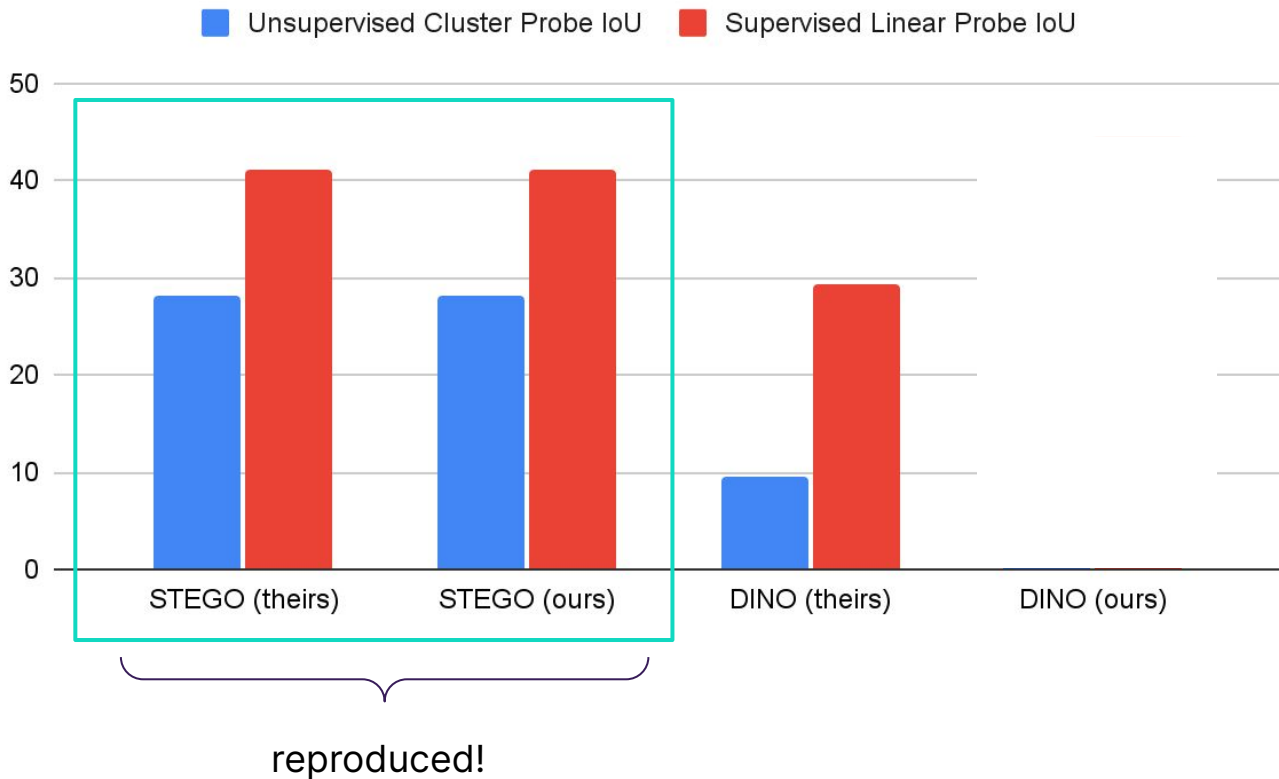- Cluster Probe maps STEGO features to ontologies using k-Means

# Reproducibility Cocostuff

**Cluster Probe** = SegHead + K-Means + Hungarian
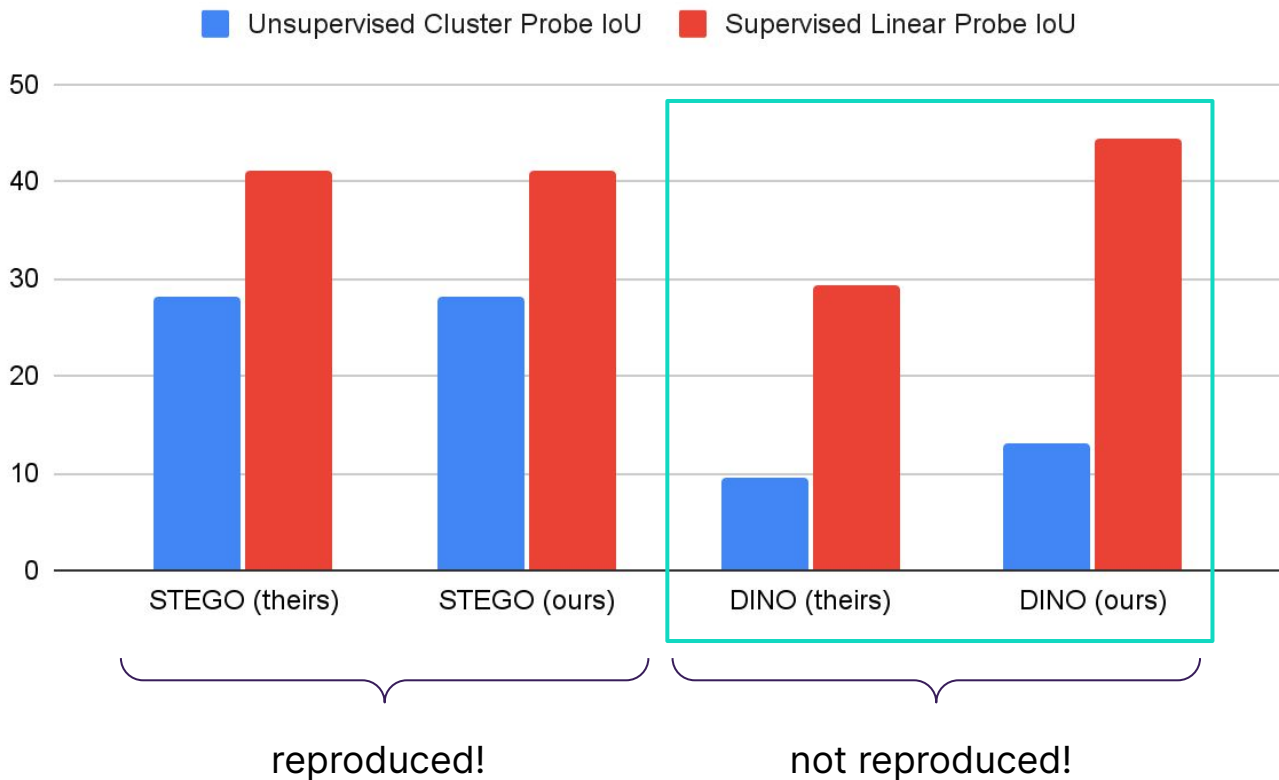**Linear Probe** = SegHead + Lin. Layer + X Entropy



Momentum

# Reproducibility Cocostuff

**Unsupervised** =  SegHead + K-Means + Hungarian
**Linear Probe** =  SegHead + Lin. Layer + X Entropy



reproduced!

# Reproducibility Cocostuff

**Unsupervised** =   SegHead + K-Means + Hungarian
**Linear Probe** =   SegHead + Lin. Layer + X Entropy



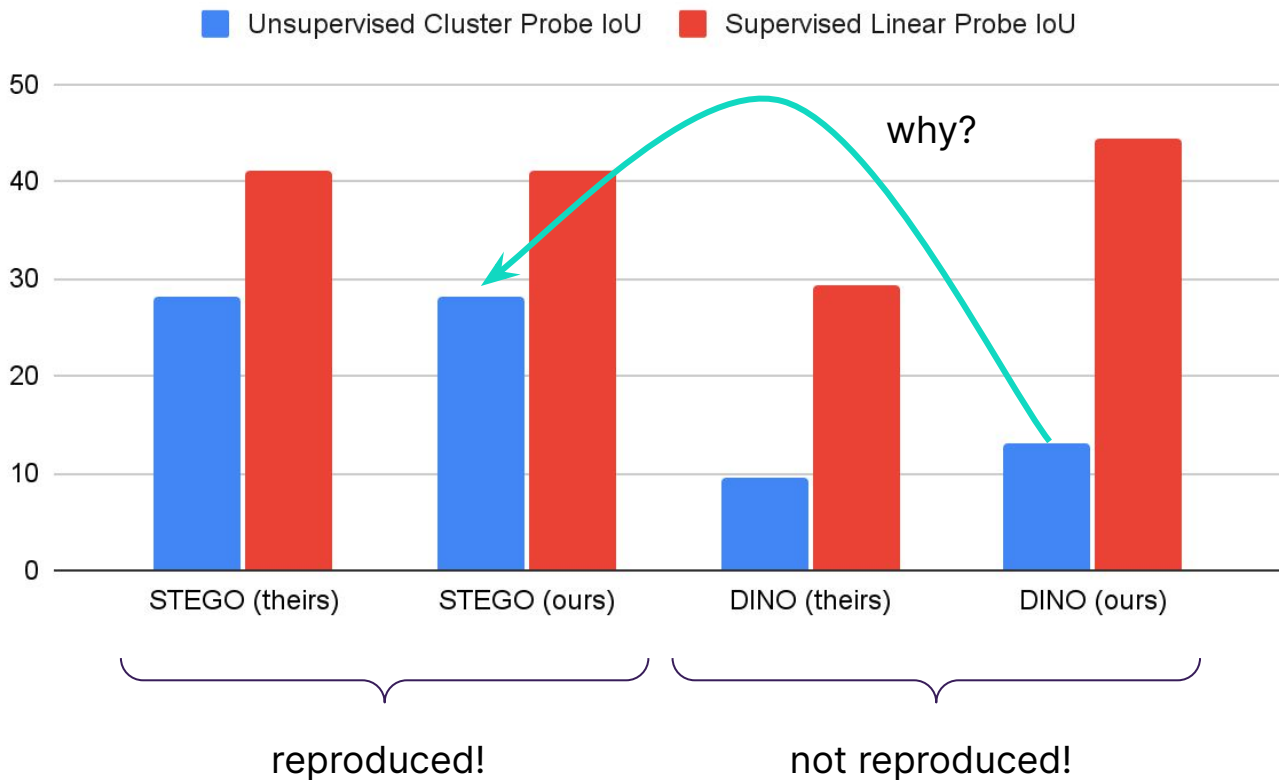■ Unsupervised Cluster Probe IoU   ■ Supervised Linear Probe IoU

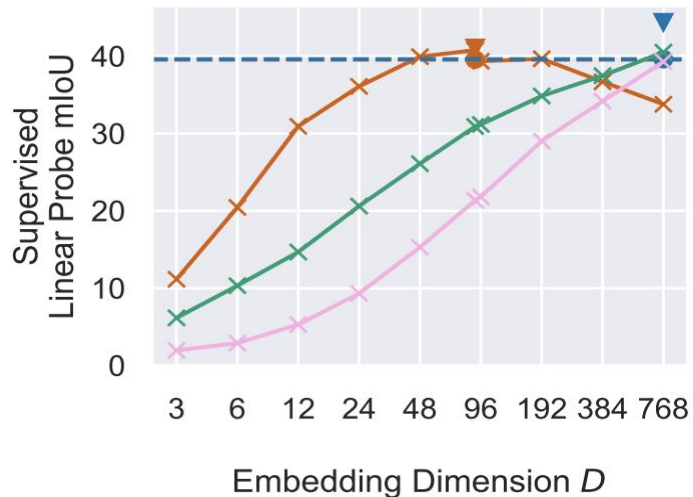reproduced!     not reproduced!

# Reproducibility Cocostuff

**Unsupervised** = SegHead + K-Means + Hungarian
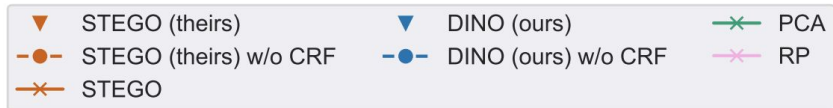**Linear Probe** = SegHead + Lin. Layer + X Entropy

# STEGO's Working Mechanisms



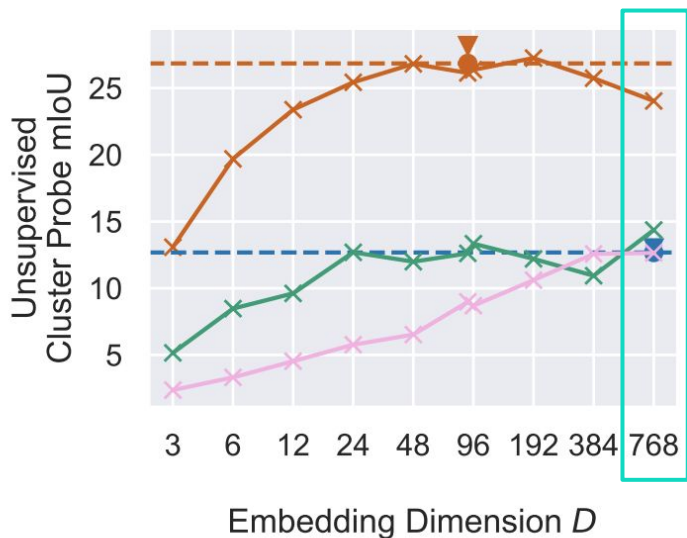## Working Mechanism 1:

- STEGO is a dimensionality reduction technique
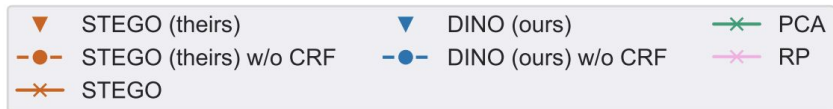- k-Means converges better in fewer dimensions

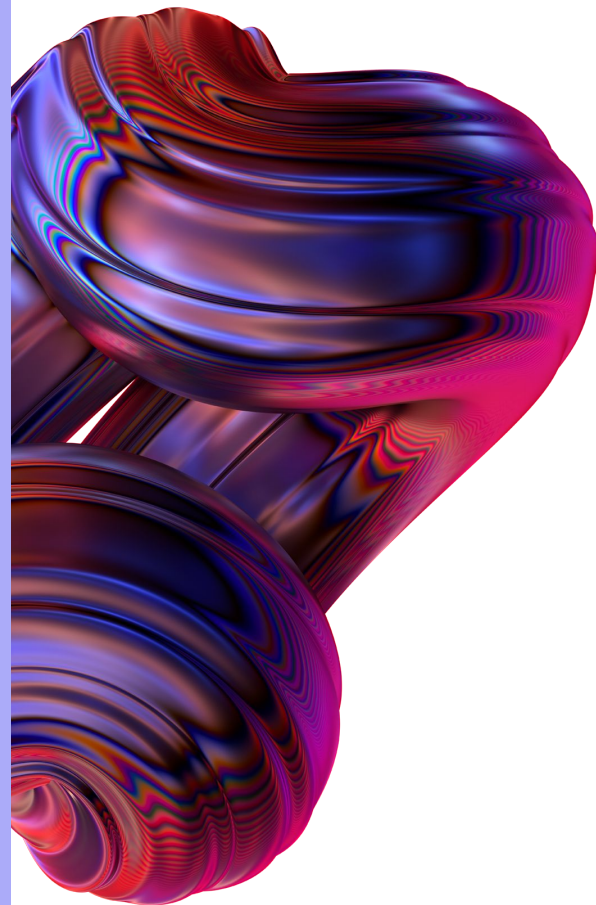# STEGO's Working Mechanisms



## Working Mechanism 2:

- Segmentation head output forms more distinct clusters

# CLIPPO: Image-and-Language Understanding from Pixels Only

Michael Tschannen, Basil Mustafa, Neil Houlsby
Google Research, Brain Team, Zürich

# CLIP-Pixels Only (CLIPPO) - Key Idea



- CLIP (Radford et al. 2021) trains separate image and text encoder

# CLIPPO - Results

- CLIPPO approaches BERT performance on GLUE benchmark
- "CLIPPO performs similarly to CLIP-style models (within 1-2%) on the main tasks CLIP was designed for - image classification and text/image retrieval"
- Good results on VQA despite never trained on that



VQAv2 dataset: Classifying CLIPPO feats

# CLIPPO - Modality Gap

CLIP* - (gap: 0.731)    CLIPPO - (gap: 0.600)    CLIPPO 25%C4 - (gap: 0.099)

**CLIPPO**

Contrastive

Transformer

CONV

Pre-training on text-text pairs with C4 (Colossal Clean Crawled Corpus) reduces modality gap

23

# CLIPPO - Typographic Attacks



| NO LABEL | | LABELED "IPOD" | | LABELED "LIBRARY" | |
|---|---|---|---|---|---|
| Granny Smith | 85.61% | Granny Smith | 0.13% | Granny Smith | 1.14% |
| iPod | 0.42% | iPod | 99.68% | iPod | 0.08% |
| library | 0% | library | 0% | library | 90.53% |
| pizza | 0% | pizza | 0% | pizza | 0% |
| rifle | 0% | rifle | 0% | rifle | 0% |
| toaster | 0% | toaster | 0% | toaster | 0% |

Source: https://distill.pub/2021/multimodal-neurons/

**Typographic attack**: "the tendency of CLIP-style models to zero-shot classify an image according to adversarially injected scene text unrelated to the scene"

**CLIPPO Result**: "All models are largely able to ignore the typographic attack, and the CLIPPO models are on par with or better than the counterparts relying on a tokenizer."

# OpenScene: 3D Scene Understanding with Open Vocabularies

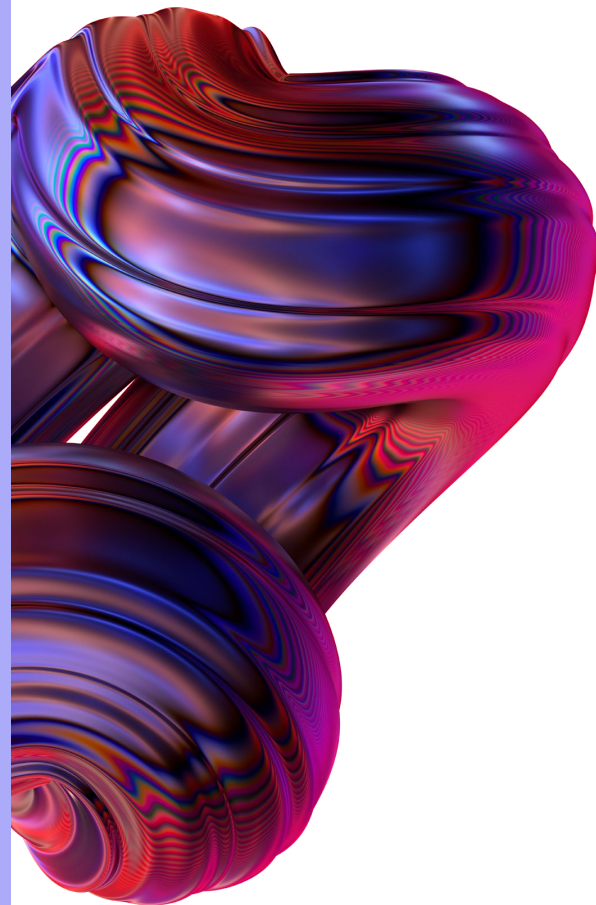Songyou Peng[1,2,3]     Kyle Genova[1]     Chiyu "Max" Jiang[4]     Andrea Tagliasacchi[1,5]

Marc Pollefeys[2]     Thomas Funkhouser[1]

[1] Google Research     [2] ETH Zurich     [3] MPI for Intelligent Systems, Tübingen     [4] Waymo LLC     [5] Simon Fraser University

pengsongyou.github.io/openscene

# Traditional (3D) Semantic Segmentation



Input 3D Geometry

Traditional Semantic Segmentation

Only train and test on a few common classes

Legend: wall, floor, cabinet, bed, chair, sofa, table, door, window, counter, curtain, toilet, sink, bathtub, other, unlabeled

# OpenScene - Key Idea



3D Geometry

CLIP Text Features
(visualize with T-SNE)

RGB Images

1. Co-embed 3D text-image features
2. Reason about properties of 3D points via cosine-similarity

# OpenScene - Demo

**Visual Programming: Compositional visual reasoning without training**

Tanmay Gupta, Aniruddha Kembhavi
PRIOR @ Allen Institute for AI
https://prior.allenai.org/projects/visprog

CVPR Award Candidate

# VisProg



- VisProg - a framework that builds CV pipelines from natural language
- "uses the in-context learning ability of GPT3 to generate python programs"
- Each line invokes functions s.a. CV models, openCV or PIL routines, …

# HandsOff: Labeled Dataset Generation With No Additional Human Annotations

Austin Xu[*]
Georgia Institute of Technology

Mariya I. Vasileva
Amazon AWS

Achal Dave[†]
Toyota Research Institute

Arjun Seshadri
Amazon Style

CVPR Highlight

# HandsOff - Key Idea



Figure 1. The HandsOff framework uses a small number of existing labeled images and a generative model to produce **infinitely** many labeled images.

- Trained on less than 50 labeled images
- GAN inversion for dataset generation

# HandsOff - GAN Inversion 101



Latent space

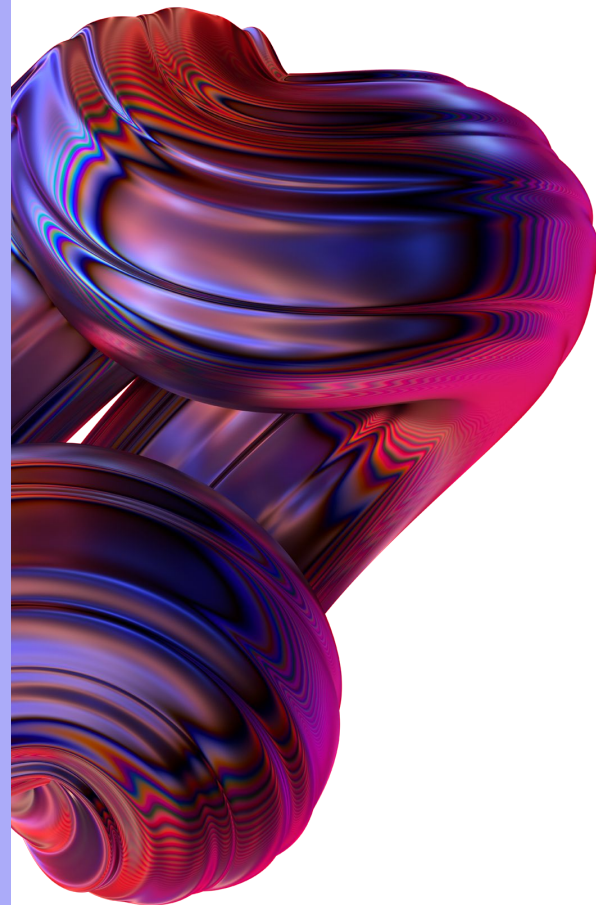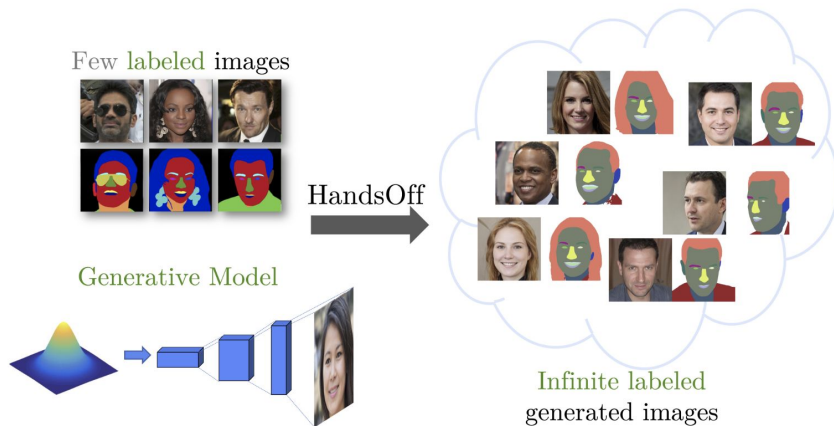Sampling & Generation

Fake Image

$x = G(z)$, $z \sim N(0, 1)$

Inversion

Real Image

(a) invert real image into latent space

$z^* = \arg\min_z (G(z), x)$

(b) manipulate the inverted image in the latent space

Reconstruction& Manipulation

$x = G(z^*)$    $x = G(z^*+n_1)$    $x = G(z^*+n_2)$

Decrease age     Add smile

**Figure 1:** GAN inversion overview

$x^{real}$    $E$    $z^*$    $G$    $x^{rec} \leftarrow x^{real}$

**Figure 2:** Invert GAN with encoder E, trained by min. rec. los

- "**GAN inversion** aims to invert a given image back into the latent space of a pretrained GAN model so that the image can be faithfully reconstructed from the inverted code by the generator"

# HandsOff - Details

grabs intermediate layers of StyleGAN2 generator, up-samples them, forms pixel-wise features

MLP head maps pixel-wise feature to label



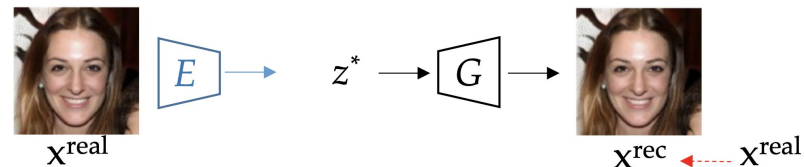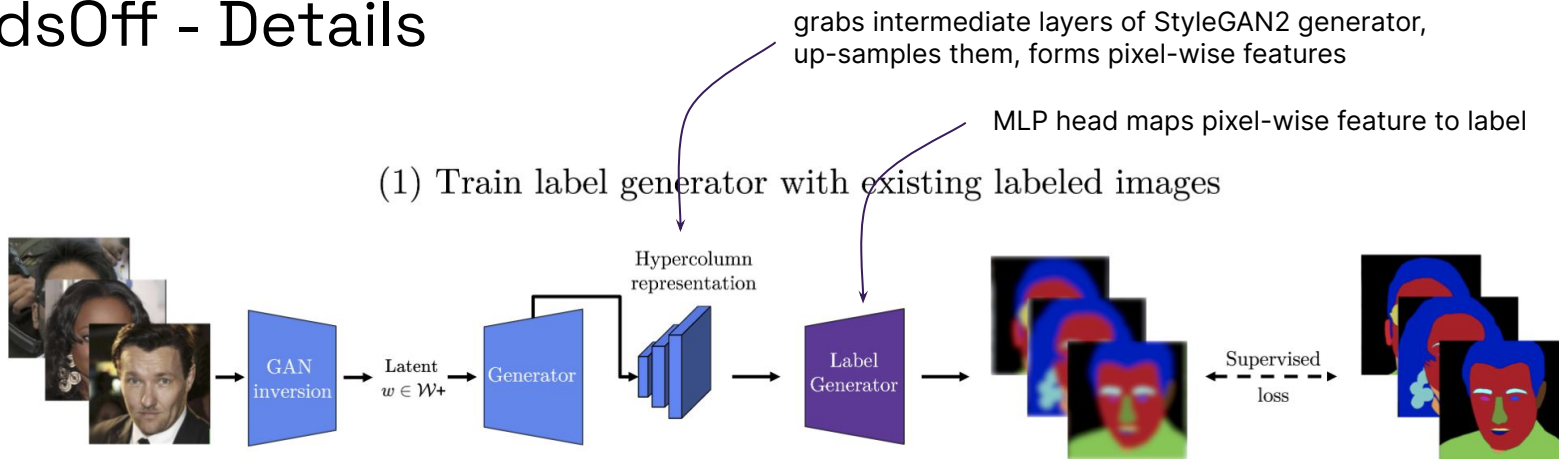(1) Train label generator with existing labeled images

GAN inversion → Latent $w \in \mathcal{W}+$ → Generator → Hypercolumn representation → Label Generator → Supervised loss

(2) Generate images and corresponding labels

Latent $w \in \mathcal{W}$ → Generator → Label Generator

# HandsOff - Results

# HandsOff - Long Tail Improvement



**Figure 1:** Long-tailed data distrubtions.



**Figure 2:** Improved Jensen-Shannon divergence and mask quality with more synthetic training data.

Source Graph:
https://www.marksayson.com/blog/advances-in-computer-vision-and-chasing-long-tail/

# ImageBind: One Embedding Space To Bind Them All
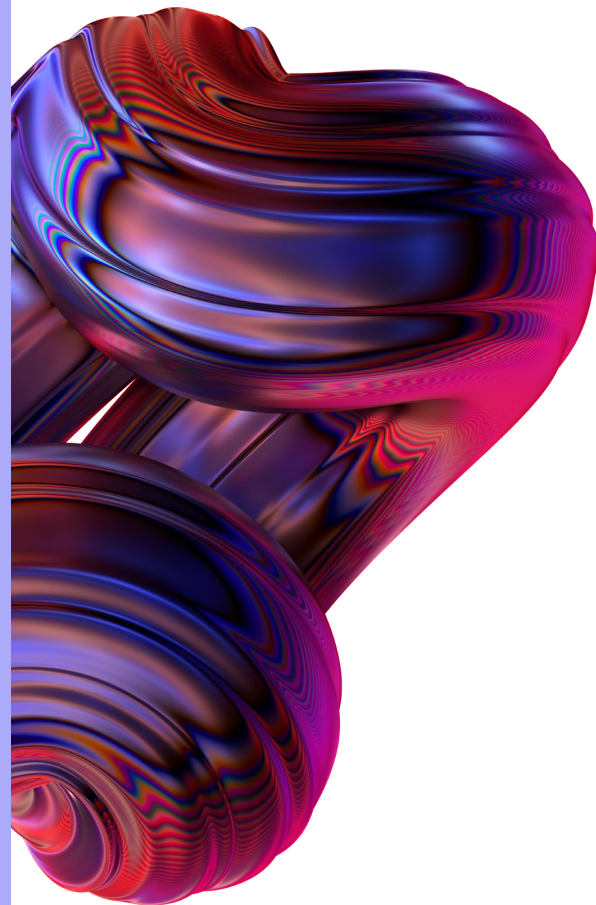
Rohit Girdhar*      Alaaeldin El-Nouby*      Zhuang Liu      Mannat Singh
Kalyan Vasudev Alwala      Armand Joulin      Ishan Misra*
FAIR, Meta AI

https://facebookresearch.github.io/ImageBind

CVPR Highlight

# ImageBind - Key Idea



**Figure 2. IMAGEBIND overview.** Different modalities occur naturally aligned in different data sources, for instance images+text and video+audio in web data, depth or thermal information with images, IMU data in videos captured with egocentric cameras, *etc*. IMAGEBIND links all these modalities in a common embedding space, enabling new emergent alignments and capabilities.

- Goal: multimodal representation learning (i.e. have single aligned feature space)
- But: no dataset couples modalities s.a. Vision, Audio, IMU, Depth, Thermal, … → self-supervision
- Idea: contrastive learning on (I, M) pairs, where I=image and M=some other modality

# ImageBind - Emergent Properties

## Cross-modal retrieval

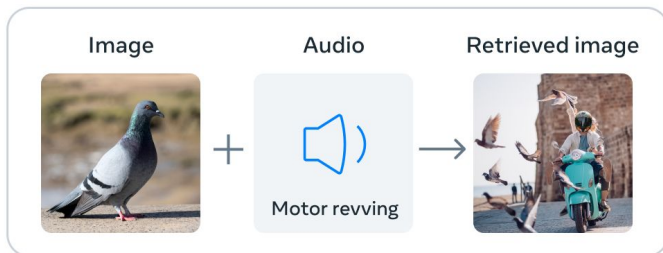| Audio | Image & Video | Depth | Text |
|---|---|---|---|
| Train horn | | | "A train pulls into a busy station"<br><br>"Wind blows as a train moves through a grassy landscape"<br><br>"People sip coffee in the dining car" |

## Embedding-space arithmetic

Image + Audio (Motor revving) → Retrieved image

## Audio to image generation

Audio (Penguin calls) → Generated image

now you can use diffusion model (DALLE-2) as image generator from audio!

See demo at: https://imagebind.metademolab.com/

# ImageBind - Emergent Properties



**Figure 5. Object detection with audio queries.** Simply replacing Detic [88]'s CLIP-based 'class' embeddings with our audio embeddings leads to an object detector promptable with audio. This requires no re-training of any model.

- ImageBind is initialized with CLIP
- Detic = pre-trained text-based detection module uses CLIP embeddings
- Idea: replace Detic's text embeddings with audio embeddings

**InstructPix2Pix: Learning to Follow Image Editing Instructions**

Tim Brooks*      Aleksander Holynski*      Alexei A. Efros

University of California, Berkeley

CVPR Highlight

# InstructPix2Pix - Goal



Figure 1. Given **an image** and **an instruction** for how to edit that image, our model performs the appropriate edit. Our model does not require full descriptions for the input or output image, and edits images in the forward pass without per-example inversion or fine-tuning.

# InstructPix2Pix - Key Idea
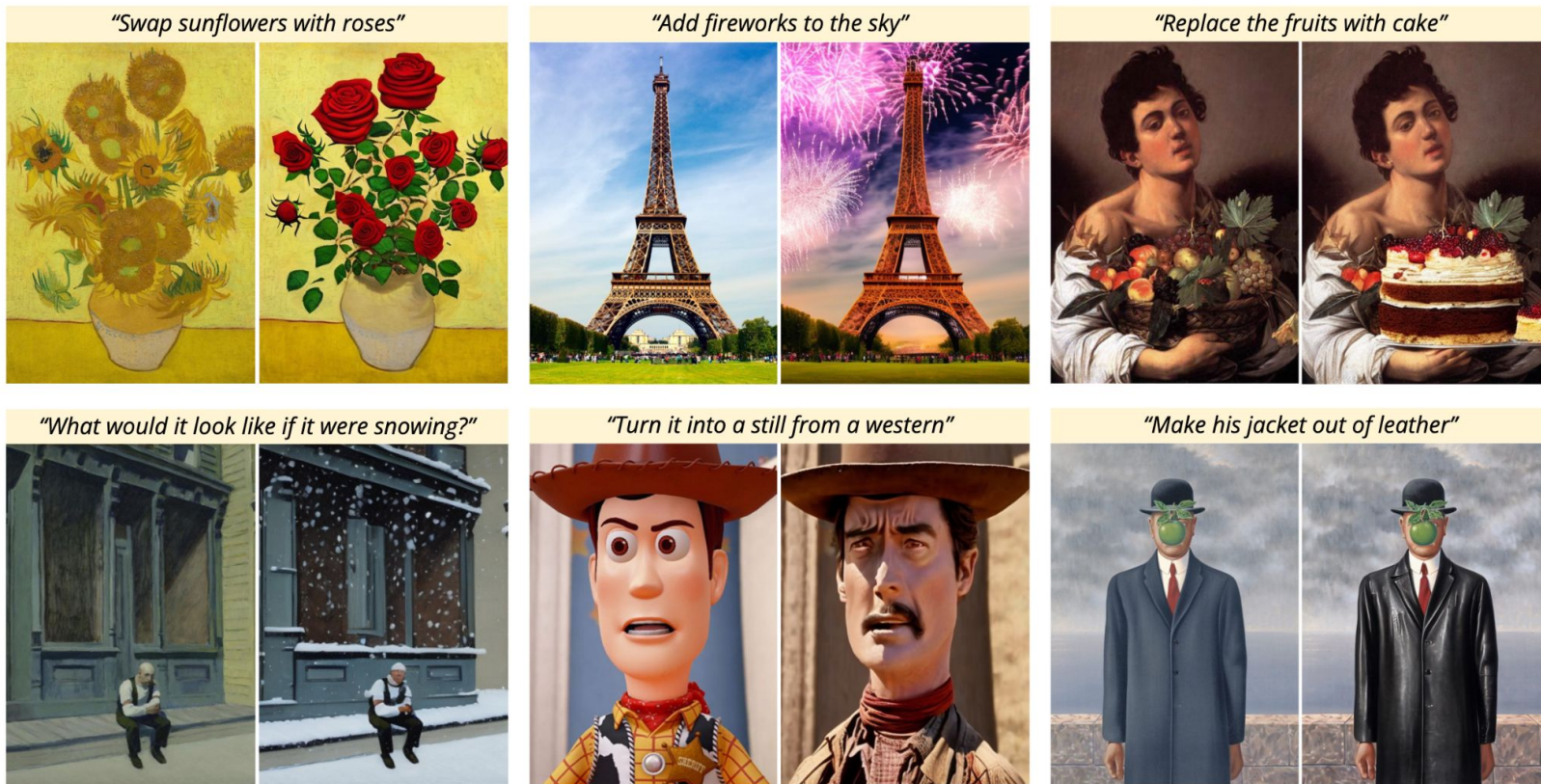


**Training Data Generation**

(a) Generate text edits:

Input Caption: *"photograph of a girl riding a horse"* → GPT-3 → Instruction: *"have her ride a dragon"*
Edited Caption: *"photograph of a girl riding a dragon"*

(b) Generate paired images:

Input Caption: *"photograph of a girl riding a horse"*
Edited Caption: *"photograph of a girl riding a dragon"* → Stable Diffusion + Prompt2Prompt

(c) Generated training examples:

*"convert to brick"*    *"Color the cars pink"*    *"Make it lit by fireworks"*    *"have her ride a dragon"*    ...

**Instruction-following Diffusion Model**

(d) Inference on real images:

*"turn her into a snake lady"*

InstructPix2Pix

- First generate 450k synthetic training samples
- Then supervised fine-tuning of pre-trained diffusion model conditioned by image
- Zero-shot generalization to real images
- But: performance is bottlenecked by models generating dataset

43

# InstructPix2Pix - More Results
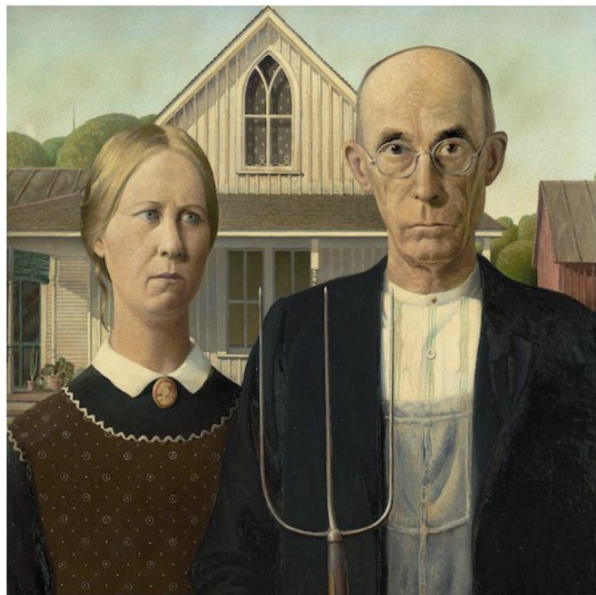


Input        *"Add boats on the water"*        *"Replace the mountains with a city skyline"*

Figure 17. A landscape photograph shown with different contextual edits. Note that isolated changes also bring along accompanying contextual effects: the addition of boats also adds wind ripples in the water, and the added city skyline is reflected on the lake.

# InstructPix2Pix - Inherited Biases



Input    "Make them look like flight attendants"    "Make them look like doctors"

Figure 14. Our method reflects biases from the data and models it is based upon, such as correlations between profession and gender.

# InstructPix2Pix - More Reading

**DreamBooth: Fine Tuning Text-to-Image Diffusion Models for Subject-Driven Generation**

Nataniel Ruiz[*,1,2]    Yuanzhen Li[1]    Varun Jampani[1]
Yael Pritch[1]    Michael Rubinstein[1]    Kfir Aberman[1]
[1] Google Research    [2] Boston University
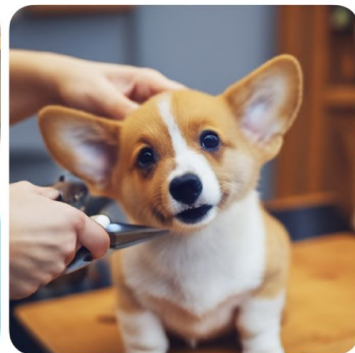
Input images

swimming  sleeping
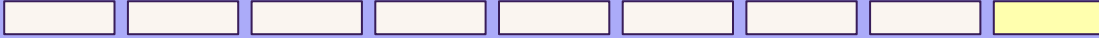
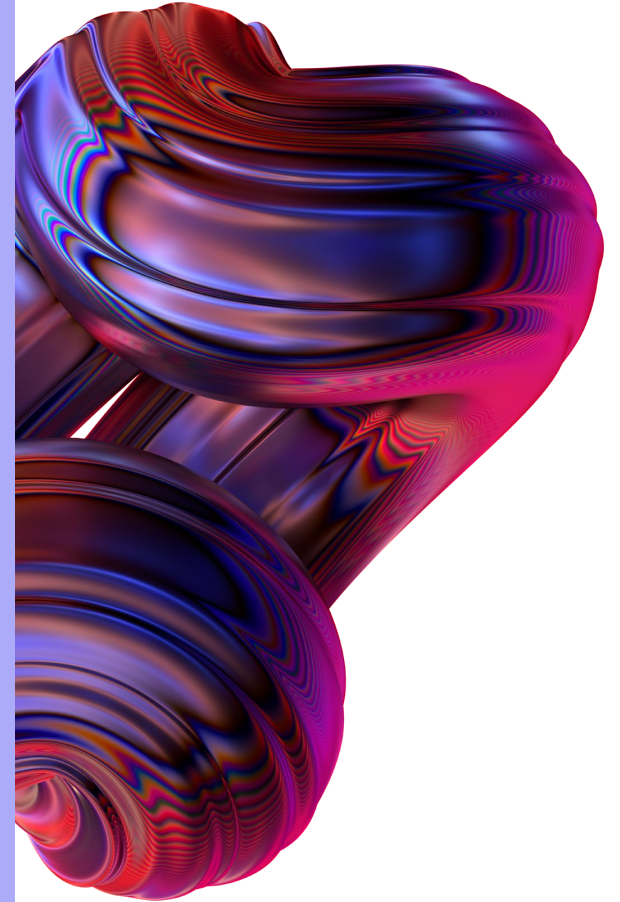in the Acropolis  in a doghouse  in a bucket  getting a haircut

Figure 1. With just a few images (typically 3-5) of a subject (left), *DreamBooth*—our AI-powered photo booth—can generate a myriad of images of the subject in different contexts (right), using the guidance of a text prompt. The results exhibit natural interactions with the environment, as well as novel articulations and variation in lighting conditions, all while maintaining high fidelity to the key visual features of the subject.
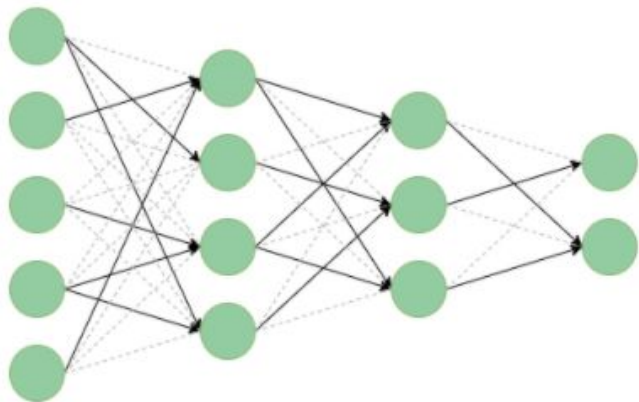
# Integral Neural Networks

Kirill Solodskikh*[†]    Azim Kurbanov*[†]    Ruslan Aydarkhanov[†]

Irina Zhelavskaya    Yury Parfenov    Dehua Song    Stamatios Lefkimmiatis
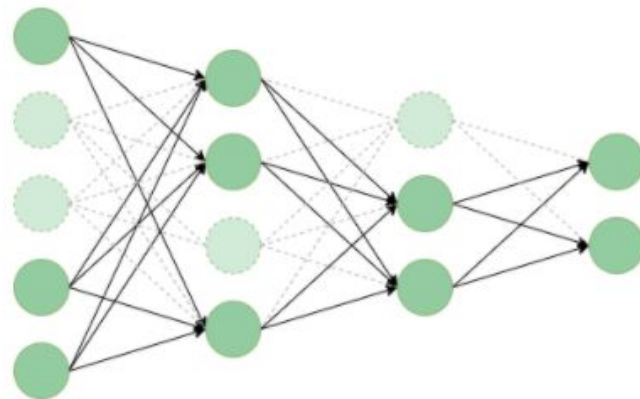
Huawei Noah's Ark Lab

CVPR Award Candidate

# Traditional Pruning



Unstructured Pruning
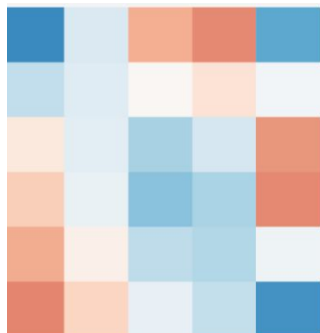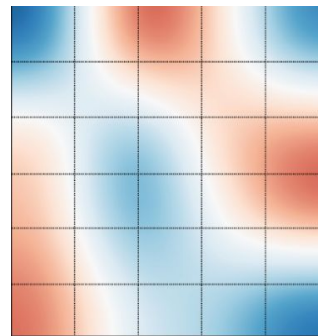
Structured Pruning

aka "weight pruning"

aka "node pruning"

# Integral Neural Networks - Key Idea



**Traditional weight tensor**
Discrete weights



**Smooth weight representation**
Sample weights at grid
Grid res. influences compression

# Integral Neural Networks

| | **Regular NN** | **Integral NN** |
|---|---|---|
| |  |  |
| **Weights** | Discrete multi-dim. tensors | Smooth multi-dim. functions |
| **Computation** | Discrete transformations of inputs | Continuous integration operations<br>Can be discretized at inference |
| **Fine-tuning** | Usually necessary after pruning | Not necessary |
| **Deployment** | Fixed model size after pruning | Resize model on-the-fly (e.g. on edge device) |

# Integral Neural Networks



Figure 1. Visualization of different channels selection methods without fine-tuning compared with our proposed integral neural networks. a) ResNet-18 on Cifar10. b) NIN architecture on Cifar10. c)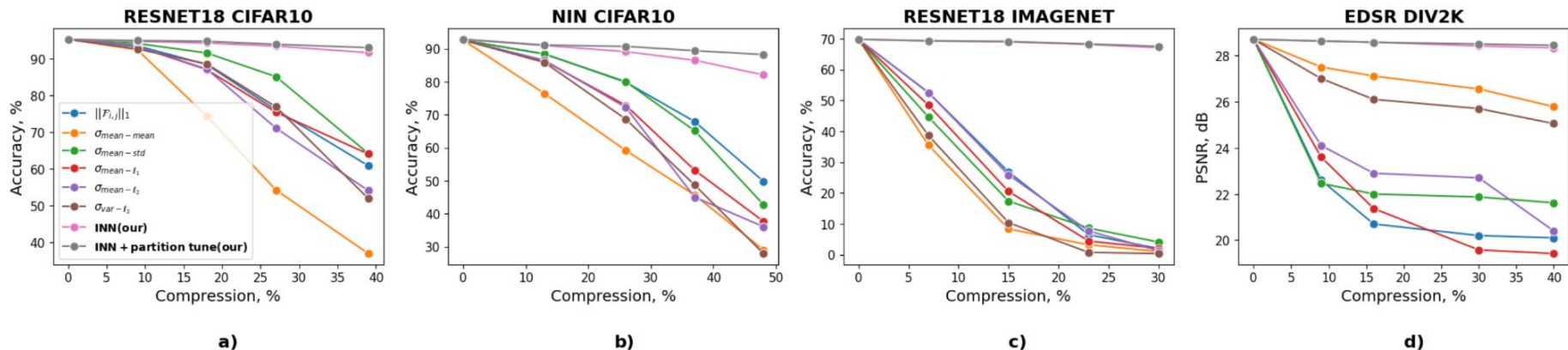 ResNet-18 on ImageNet. d) 4x EDSR on Div2k validation set. By compression we denote the percentage of deleted parameters.

# Compressing pre-trained nets



**Matrix of the original discrete weights** — $c^{out}$, $c^{in}$

**Permutation algorithm** →

**Permuted discrete weights** — $c^{out}$, $c^{in}$

**Smooth interpolation** →

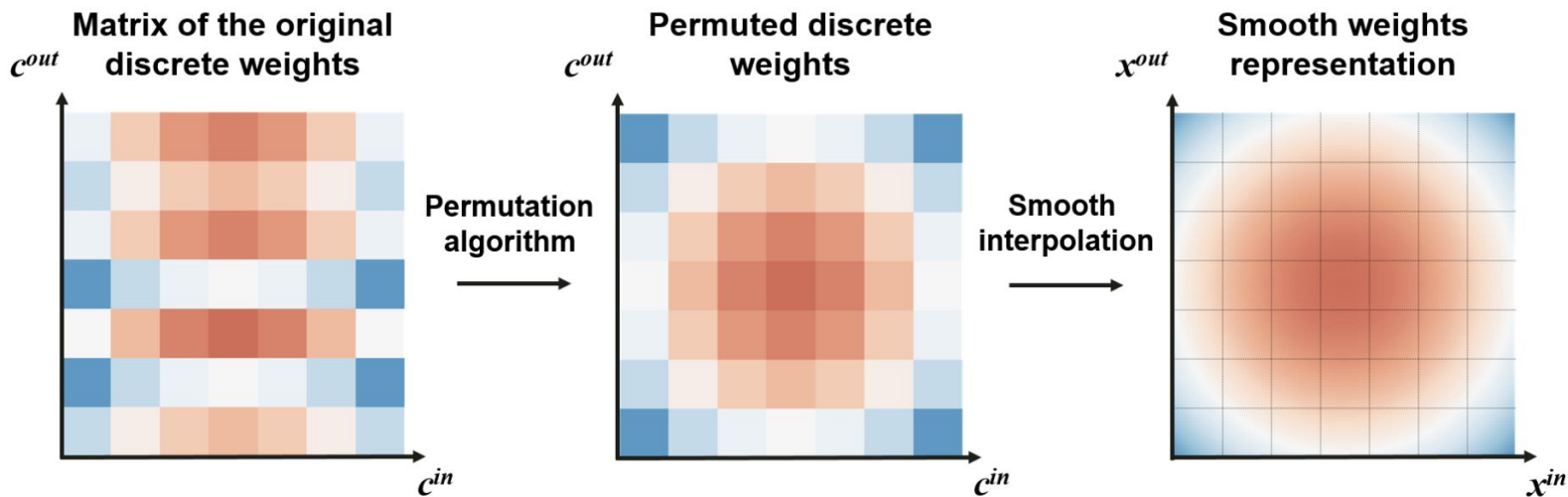**Smooth weights representation** — $x^{out}$, $x^{in}$

Figure 5. Toy example illustrating the permutation of filters in a discrete weight tensor in order to obtain a smoother structure.

CVPR
Impressions

Papers

**Food for
Thought**

**Embodied Foundation Models**

Robotics and CV

AI and "The Hype"

# Foundation Models as Motion Planners (1/2)



Driess et al. 2023 "PaLM-E: An Embodied Multimodal Language Model"

# Foundation Models as Motion Planners (2/2)

- VQA Model by W

# "World Models" - GAIA-1 by Wayve

- World model = a generative model that predicts what happens next conditioned on an action
- Autoregressive model trained on Wayve's large unlabeled dataset

# "World Models" - GAIA-1 by Wayve



- Autonomous driving may be the first example of where we see embodied AI working

| CVPR Impressions | Papers | **Food for Thought** |
|:---:|:---:|:---:|

Embodied Foundation Models    **Robotics and CV**    AI and "The Hype"

# Robotics and Computer Vision

- "Robotics is the next big thing"
- Jitendra Malik
  - ~ "Vision has no use by its own. It needs to guide action."
  - ~ "Robotics is 20 years behind computer vision"
  - ~ "Navigation and locomotion are close to being solved."
  - ~ "Manipulation is far from being solved" → Why?
    - Control struggles with making and breaking of contact
    - RL struggles with inaccurate simulations and sim to real gap
    - Lack of dexterous multi-fingered hands
  - Urged the CV community to venture into manipulation
- Differences CV and robotics
  - no standardized benchmarks
  - no large datasets
  - sim to real gap
  - hardware experiments are essential but take long
- Particular hot topics: visual pre-training for robotics, object representations



Train in Simulation

Qi et al. CoRL 2022

CVPR
Impressions

Papers

**Food for
Thought**

Embodied Foundation Models

Robotics and CV

**AI and "The Hype"**

# AI and "The Hype"
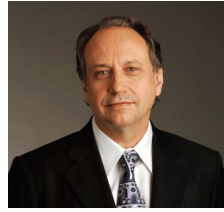
- Rodney Brooks, also see "The Seven Deadly Sins of AI Predictions" blog
    - Roy Amara (1925 – 2007): "We tend to overestimate the effect of a technology in the short run and underestimate the effect in the long run."
    - Example: Fear that computers will replace librarians, librarians kept on going for 40 years until eventually largely impacted by the internet and mobile devices.



IJCAI 1979, Tokyo

AN AUTOMOBILE WITH ARTIFICIAL INTELLIGENCE

Sadayuki Tsugawa, Teruo Yatabe, Takeshi Hirose and Shuntetsu Matsumoto
Automobile Division
Mechanical Engineering Laboratory
5-12-2 Fujimi-cho, Higashimurayama
Tokyo 189 Japan

This paper describes an automobile with artificial intelligence, which consists of a road pattern recognition unit and a problem solving unit. The vehicle is completely autonomous and can be driven without a human driver. The road pattern recognition unit involving a pair of TV cameras and a processing unit identifies obstacles in front of the vehicle and outputs data regarding to the locations of the obstacles. The problem solving unit is a microcomputer system and determines control optimal to the environment around the vehicle based on the data. The algorithm employed in it is a table-look-up method, in which the location of the optimal control is addressed in the table by key words generated from the data. The table was heuristically made by means of digital simulation. The vehicle was successfully driven under various road environments at the speed within 30



FORECASTS: http://www.driverless-future.com/?page_id=384 March 27, 2017

NVIDIA to introduce level-4 enabling system by 2018 (2017)
NuTonomy to provide self-driving taxi services in Singapore by 2018, expand to 10 cities around world by 2020 (2016)
Delphi and MobilEye to provide off-the-shelf self-driving system by 2019 (2016)
Ford CEO announces fully autonomous vehicles for mobility services by 2021 (2016) ←
Volkswagen expects first self driving cars on the market by 2019 (2016)
GM: Autonomous cars could be deployed by 2020 or sooner (2016) ←
BMW to launch autonomous iNext in 2021 (2016) ←
Ford's head of product development: autonomous vehicle on the market by 2020 (2016) ←
Baidu's Chief Scientist expects large number of self-driving cars on the road by 2019 (2016)
First autonomous Toyota to be available in 2020 (2015) ←
Elon Musk now expects first fully autonomous Tesla by 2018, approved by 2021 (2015)
US Sec Trans: Driverless cars will be in use all over the world by 2025 (2015)
Uber fleet to be driverless by 2030 (2015) ←
Ford CEO expects fully autonomous cars by 2020 (2015) ←
Next generation Audi A8 capable of fully autonomous driving in 2017 (2014)
Jaguar and Land-Rover to provide fully autonomous cars by 2024 says Director of Research and Technology (2014)
Fully autonomous vehicles could be ready by 2025, predicts Daimler chairman (2014) ←
Nissan to provide fully autonomous vehicles by 2020 (2013) ←
Truly autonomous cars to populate roads by 2028-2032 estimates insurance think tank executive (2013)
Continental to make fully autonomous driving a reality by 2025 (2012)

"Don't be the best, be the only!"

# Thanks!

Alexander Koenig
alexander.koenig@merantix.com