Research Report

Self-supervised Point Representations for Spatial Cognition

Xiaoyang Wu June 9, 2025





DINOv2 is changing how we research

We were initially impressed by its **semantically meaningful visualization** of representations We found it to be super robust in many tasks, even with a **simple linear probing**.







Retrieval





Sparse matching

Dense matching





DINOv2 is changing how we research

If we process images, we **absolutely consider** using it, and it rarely disappoints us! This shapes trust and reliability in our minds!





PixelSplat (Reconstruction)

> Cambrian-1 (Multimodal)







DINO Series (Image SSL)



Depth Anything (Perception)



AnyDoor (Generation)







Two key factors contribute to the powerful capabilities of representation



DINO Series (Image SSL) The magic of data: The foundation of the intelligence

142M High Quality Images

2 web	Detairt / Split	- Interaction	Retrieval	Batricynd	Final
clanification.	ImageNet-22k /	14,197,086	an in		14,197,086
classification	ImageNet-IIk / -	14.197,086	animple	36,768,344	56,768,544
(lassification)	ImageNet-1k / train	1,291,147	sample	41,997,344	81,397,344
Bas-grained chevel.	Culturk 101 / train	3,030	cluster	2,430,000	1,000,000
for-grained slassif.	CUB-200-2011 / train	5,994	(limiter	1,380,000	1,800,000
fine-grained classif.	DTD / train1	1,880	chaster	1.540.000	1,000,000
Riv-gratied classif.	PGVC-Astends / train	3,334	chetre	1,179,008	1,000,000
fan-grained classif.	Flowers-302 / train	1,630	cluster	1,060,000	1,800,000
first-grained classif.	Food-101 / train	75,790	choise	21,670,000	1,000,000
fine-grained classif	Oxford-IIIT Pet / trained	3,680	chaster	2,750,000	1,000,000
fine-grained elassif.	Standard Care / train	8,546	chains	7,239,008	1,000,000
for-grained classif.	SUNDE / undel	19,850	cheetere	18,300,000	1,000,060
fae-grained classif.	Pascal VOC 2007 / train	2,501	rbeter	3,0109,0000	1,000,000
regneralation.	ADEDIK / traile	20,210	chater	20,730,000	1,800,000
anguamitation.	Cityarapes / train	2,975	chater	1,290,000	1,800,000
erginetation.	Pheral VOC 2013 (seg.) / training	1,604	distort	15,140,000	1,000,000
depth estimation	Mugillary SLS / train	1,434,292	ad in		1,414,242
depth rotinuation	KITTE / trade (Eigen)	23,154	chaire .	3,700,009	1,000,000
depth estimation.	NYU Depth V2 / train	24,231	cluster	TE.A50.000	1,000,000
depth estimation	SUN BOB-D / train	4,829	diater	4,375,008	1,800,000
retrieval	Google Landmarks v2 / train (clean)	1,540,470	and its		1,345,470
set/inval	Google Landmarks v2 / train (ritors)	1,540,470	autopla	8,331,680	8,321,880
average of the second s	AnoterTine / new	1,238	elimites.	960,000	960,000
estrival	AmsterTime / old	1,231	chaster	830.000	830.000
ant/anal	Mick / train	397,121	chaiter	62,805,000	1,000,000
petricoul.	Bevisiting Oxford / hum	4,990	(buter)	3,680,000	1,000,000
retrieval	Revhilling Parls / base	4,303	cluster	3.660.000	1,000,000
					141,100,3m

The magic of self-supervised learning: The potential of leverage "infinite" data







There are so many spatial tasks we care about!







Xiaoyang's Research Report



大 學



Background What's the DINO time of 3D?

It leverages the magic of data and learning; It (target at) unify and benefit all 3D tasks;





大學

Background We can absolutely solve them using images, but ...





While sometimes it is good, several restrictions remain:

- It requires significant inference time to process a scene video with thousands of frames.
- It lacks geometric awareness, as it is not natively encoded in 3D space.





We believe the point cloud is an ideal data structure for spatial cognition.

Why point cloud?

- Point cloud is everywhere and easy to acquire





Laser Scanner









We believe the point cloud is an ideal data structure for spatial cognition.

Why point cloud?

Now we can also get 3D/4D point cloud from dynamic videos



CUT3R





Let's see what kind of representations we can acquire from these sparse point clouds.



Similarity Heatmap







Let's see what kind of representations we can acquire from these sparse point clouds.

In our latest research (will be released in July)



Semantic Segmentation			ScanNet Val [15]			
Method	Туре	Encoder	mIoU	mAcc	allAcc	
DINOv2 [26]	2D Image SSL	ViT-G	63.09	75.50	82.42	
Sonata [47]	3D Point SSL	PTv3-B	72.52	83.11	89.74	
Sonata×DINOv2	3D SSL×2D SSL	Both	75.91	85.36	91.25	
Concerto (ours)	2D-3D Joint SSL	PTv3-B	77.32	86.58	91.74	

Linear Probing

PCA Visualization

Xiaoyang's Research Report



大學



Lifting and encoding images from scene videos

In our latest research (will be released in July)



Xiaoyang's Research Report



大 學



Probing representations into language space

In our latest research (will be released in July)



Xiaoyang's Research Report



大學



Roadmap

Excited about these results? Let's dive into the technical details!





ERSITY OF HONG KONG



Paper Point Transformer V3: Simpler, Faster, Stronger

- \Rightarrow State-of-the-art performance on over **20** downstream tasks that span both **indoor** and **outdoor** scenarios.
- \Rightarrow Expanding the receptive field from **16** to **1024** points while remaining efficient.
- \Rightarrow 3x increase in processing speed and 10x improvement in memory efficiency compared with PTv2





We were suffering from an extremely small scale of data!



 \Rightarrow The largest (2022) indoor perception dataset, ScanNet v2, only contains 1,613 scene-level point clouds.

 \Rightarrow The largest dataset only worths 1,613 images with 300×300 resolution?





Backbone designs are also overfitting to complex designs!



- ⇒ Validation loss increase during training (observed in PTv2, the issue is almost addressed after PPT and PTv3)
- \Rightarrow Why can not unify parameters for different tasks
- \Rightarrow Why backbones overfitting to complex designs

Point Prompt Training, by Xiaoyang Wu, et al.; accepted by CVPR 2024





Enlarging data scale and training scale with multi-dataset Point Prompt Training



Point Prompt Training, by Xiaoyang Wu, et al.; accepted by CVPR 2024





Philosophy PTv3: Scaling Principle for Backbone Design



Scaling Principle

- \Rightarrow Model performance is **more significantly influenced by scale** than by complex design;
- \Rightarrow We should **prioritize simplicity and efficiency over the accuracy** of certain mechanisms;
- \Rightarrow Efficiency enable scalability and further enable stronger accuracy.





Motivation Breaking the curse of permutation invariance



- \Rightarrow The unstructured nature of point cloud makes point-based methods rely on K-nearest neighbor.
- \Rightarrow However, K-nearest neighbor take up to **28%** forward latency!





Breaking the curse of permutation invariance



Point Transformer V2



Point Transformer V3

- \Rightarrow Do we really need the accurate neighbors?
 - \Rightarrow **No**, attention is adaptive to kernel shape, large kernel >> precise neighbors;
- \Rightarrow Serialized point cloud into a structured format.





Method Space-filling Curve







Method Point Cloud Serialization



Point Cloud



Serialized Point Cloud structured 1D format





Method Serialized Attention









Method Serialized Attention







One more thing

The potential for multi-modal feature extraction

- \Rightarrow Point cloud high-dimensional sparse data.
- \Rightarrow Many high-dimensional data with positional information can be view as a point cloud.
 - \Rightarrow Like image, climate and of course point cloud



Image (2D)



Climate (2.5D)



Point cloud (3D)





Background What's the DINO time of 3D?

It leverages the magic of data and learning; It (target at) unify and benefit all 3D tasks;

It is Project Sonata





大學



Paper Sonata: Self-Supervised Learning of Reliable Point Representations







What is Self-supervised Learning?

=> Make things should be same, the same.

"No plain, no gain", the more the model struggles, the better the representations it learns.
 Just like a child, the model is good at finding shortcuts.



Local Crop

Student (active learning)

Teacher (EMA Update)

Global Crop

Contrastive Learning as an example





What is Self-supervised Learning?

=> Make things should be same, the same.

=> Just like a child, the model is good at finding shortcuts.



Contrastive Learning as an example

Xiaoyang's Research Report



大學

PointContrast, by Saining Xie, et al.



Why we introduce strong augmentations and mask image/point modeling?



Core Principle: continuously increasing the difficulty of pretext tasks as long as the model continues to converge Local-Global Alignment: The "base" in the "cocktail" of SSL, setting the foundation and direction for convergence. Masked-Unmasked Alignment: The "liqueur" that intensifies the challenge of SSL, adding complexity.





Why we introduce strong augmentations and mask image/point modeling?









Brightness +40%



Brightness -40%







Saturation +20%

Contrast +40

Saturation -20%





Masked Scene Contrast, by Xiaoyang Wu, et al.



∧ Meta



Why we introduce strong augmentations and mask image/point modeling?



Xiaoyang's Research Report



大 學



Why we introduce strong augmentations and mask image/point modeling?











Local Crop

(active learning)

(EMA Update)

Global Crop

Teacher (backup): a teacher is actually a kind of "backup" and "smooth average" of a student; **Student (spearhead):** the student is more like a "spearhead", tries something challenging and risky;

With the teacher preventing the student from being misled, the student is less likely to get lost in "mission impossible" and has a greater chance to discover treasure within "impossible"





Method Why Self-distillation?



Now we lead to the DINOv2 Architecture...





What is **Point Cloud** Self-supervised Learning?

Is it just adapting Image SSL to point cloud data?





Challenge What is the Geometric Shortcut





* With 0.02% Learnable Parameters Linear Probing

Geometric shortcut **refers to** the tendency of the model to **collapse to easily accessible, low-level geometric cues**, such as normal direction or point height. Therefore, resulting low linear probing performance in downstream tasks





Challenge

How do geometric shortcuts occur



Point Transformer V2



Point Transformer V3

Evidence: The point cloud encoder can capture geometric information not only from input features but also from the operator kernel defined by coordinates.





Challenge

The Geometric Shortcut is a unique challenge to Point Cloud (Sparse Data) 🎊

Why not happen with image SSL

> Image A' and Image A are same. > Image **B** and Image **A** are different.

> Mask out all the feature

- > Shuffle the order of three images.
- > Can we distinguish which is A and A'?
- > No, we can't without any feature
- > How about add some feature back?
- > Yes, now we can!!
- >> Image SSL have to relies on feature and impossible to collapse into geometry information (coordinate)



Image A' Image A Image B shuffle shuffle Image ??? Image ??? Image ??? Image ??? Image ??? Image ???



Challenge

The Geometric Shortcut is a unique challenge to Point Cloud (Sparse Data) 🌋





Point Cloud A



Point Cloud A'

Point Cloud **X**?



Point Cloud **P**?









- > PC A' and PC A are same.
- > PC **B** and PC **A** are different.
- > Mask out all the feature

- > **Shuffle** the order of three images.
- > Can we distinguish which is A and A'?
- > Yes, we can. Easy!

- > Do we really need other feature?
- > No, no need.
- >> Point cloud SSL exist a shortcut that only rely on geometry information



Challenge The Geometric Shortcut is a unique challenge to Point Cloud (Sparse Data)



Xiaoyang's Research Report



大學



Key solutions to geometric shortcut



Core concept: Disturb on positional information to restraint collapse



Xiaoyang's Research Report



大學



Key solutions to geometric shortcut



Core concept: Disturb on positional information to restraint collapse



Good: Improve Linear probing from 20% to 60% +

Bad: Here only contains global-wise embedding, also need some local information





Key solutions to geometric shortcut



Core concept: Disturb on positional information to restraint collapse





Method Other solutions to geometric shortcut (Micro Designs)



Mask ratio scheduler & Mask size scheduler: Set a trap to against geometric shortcut



Mask size: 0.1m Mask ratio: 30%



Mask size: 0.4m Mask ratio: 70%

Xiaoyang's Research Report

Random gaussian jitter with masked point: Break local geometric information





Scaling **Data Scaling**







One more thing Comparison with DINO in 3D Space

Semantic Segmentation			ScanNet Val [15]		
Method	Туре	Encoder	mIoU	mAcc	allAcc
DINOv2 [26]	2D Image SSL	ViT-G	63.09	75.50	82.42
Sonata [47]	3D Point SSL	PTv3-B	72.52	83.11	89.74
Sonata×DINOv2	3D SSL×2D SSL	Both	75.91	85.36	91.25



香港大學 THE UNIVERSITY OF HONG KONG



Paper The power of multi-sensory learning



- When we see an image of an apple, even in grayscale, it is easy for us to **imagine** its geometric **shape**, its possible **color**, and even its unique favors

To see a World in a Grain of Sand And a Heaven in a Wildflower - William Blake





Paper The power of multi-sensory learning



Why human able to produce such powerful embedding?
 We see apple before, we touch apple before, we eat apple before. We know these multi sensory experience. We formulate a unified representation in our mind. (SSL)







Paper

Concerto: Joint 2D-3D Self-Supervised Learning Emerges Spatial Representations





大 學

港



Paper Further scaling up the training with video dataset



Xiaoyang's Research Report



港

大 學



Future Native 2D-3D Joint Self-supervised Learning

Bidirectional Joint Embedding Learning



Two module can be asynchronous





Thank you!

cept / Pointcept Public				☐ facebookresearch / sonata Public			☐ Notifications ♀ Fork 9 ☆ Star 34
⊙ Issues 264 11 Pull requests 5	5 💿 Actions 🗄 Projects 💿 Security 🗠 Insights			<> Code issues 16 Pull requests issues 16) Actions 🖽 Projects 🛈 Security 🗠 Insights		
🍹 main 🔹 🥲 1 Branch 🛇 9 Tags	Q. Go to file	↔ Code 👻	About	🐉 main 👻 🐉 1 Branch 📎 0 Tags	Q. Go to file	<> Code →	About
👔 birgerbr Fix undefined data_part in GridSample train (#461) 🚥 🧹 3636d26 - last month. 🕲 112 Commits		Pointcept: a codebase for point cloud perception research. Latest works: Sonata (CVPR'25 Highlight), PTv3	Gofinge Add batched forward demo	18c09ff - 5	days ago 🕚 8 Commits	[CVPR'25 Highlight] Official repository of Sonata: Self-Supervised Learning of	
github/workflows	Fix formatter check branch	2 years ago	(CVPR'24 Oral), PPT (CVPR'24), MSC	github	Update CONTRIBUTING.md (#21)	last week	Reliable Point Representations
Configs	Add SGIFormer	last month	(CVPR 23)	assets	initial commit	3 months ago	
🖿 libs	Add SGIFormer	last month	point-cloud pytorch 3d-vision				C Readme
pointcept	Fix undefined data_part in GridSample train (#461)	last month	Readme	demo	Add batched forward demo	5 days ago	화 Apache-2.0 license
scripts	set ulimit in launch scrints	2 months ann	4] MIT license	📄 sonata	Add batched forward demo	5 days ago	Code of conduct
	Add Instance Commutation Tester with submission summe	2 months and	Activity Custom properties	🗅 .gitignore	initial commit	3 months ago	An Activity
	Aud instance segmentation rester with submission support.	z montris ago	☆ 2.2k stars		initial commit	3 months ago	E Custom properties
	Release Point Prompt Training (PPT)	2 years ago	 20 watching 	2			公 340 stars
LICENSE	Initial commit	2 years ago	¥ 256 forks	C README.md	Fixing various typos (#22)	last week	 17 watching
README.md	Update Readme	2 months ago	Report repository	🗋 environment.yml	Add batched forward demo	5 days ago	양 9 forks
C environment.yml	Add support to wandb (#413)	2 months ago	Releases a	🗅 setup.py	initial commit	3 months ago	Report repository

pointcept/pointcept (Pre-training)



香港大學 THE UNIVERSITY OF HONG KONG

facebookresearch/sonata (Inference)

