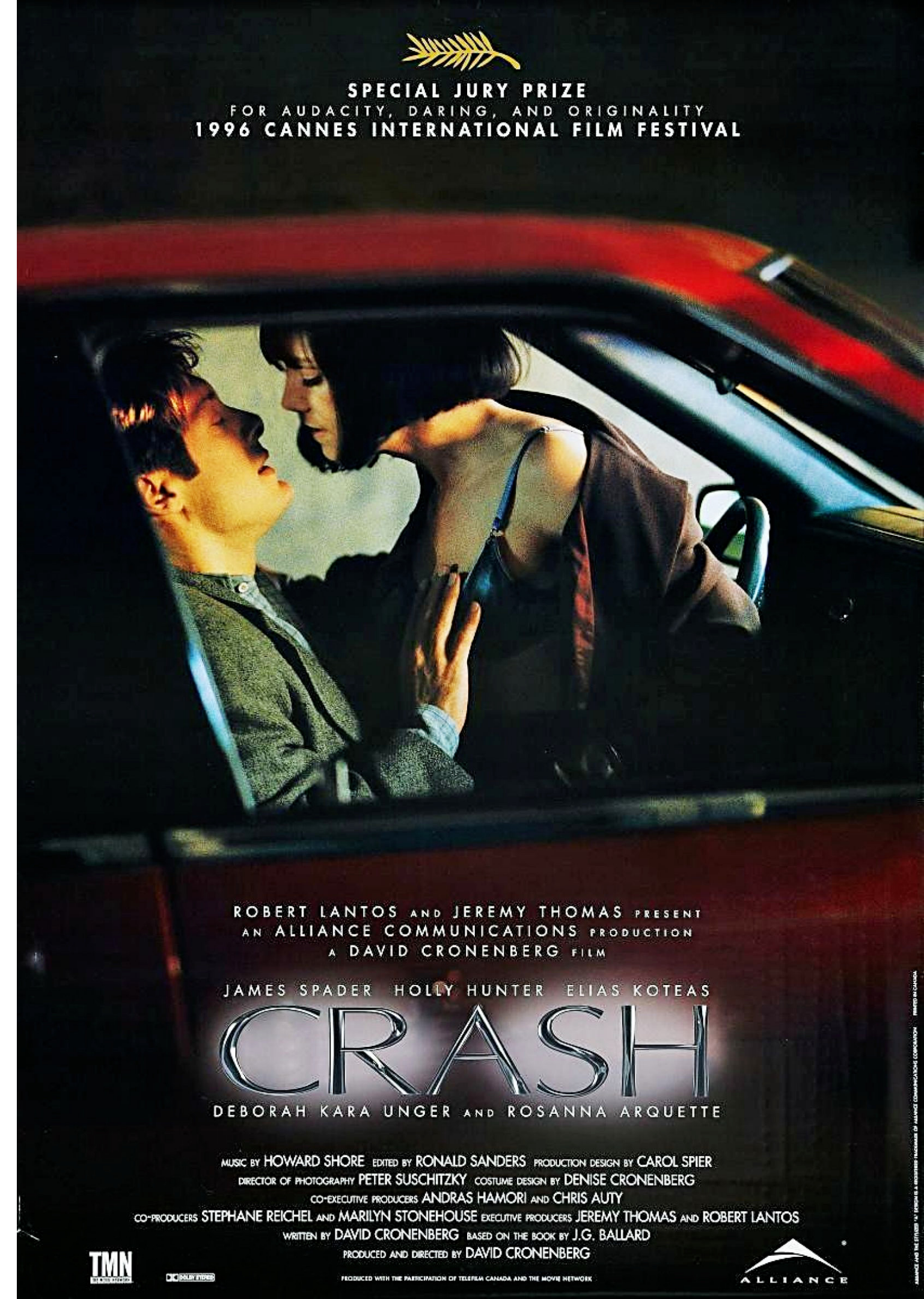


Learning Image Representations Without Manual Annotations and Scientific Applications

Piotr Bojanowski, FAIR, Meta

Manual annotations...

- Hollywood 2
 - Marcin Marszalek, Ivan Laptev, and Cordelia Schmid. "Actions in context." In *CVPR 2009*.
 - 810 + 884 videos
 - 12 actions
 - 69 Hollywood movies
- Hollywood 3 ?
 - Annotate all movies exhaustively
 - In charge of one of the movies



Scarcity of Manual Annotations

- Annotations are expensive (if high quality)
- Are ambiguous
- Class definitions are not static
- Intractable with increasing complexity of the task

Baking the Cake

- Supervised Learning is needed!
- Unsupervised learning should do the heavy lifting
- Modern success of LLMs follows this exact recipe...
- Is the vision cake ready?

“Pure” Reinforcement Learning (cherry)

- ▶ The machine predicts a scalar reward given once in a while.
- ▶ **A few bits for some samples**

Supervised Learning (icing)

- ▶ The machine predicts a category or a few numbers for each input
- ▶ Predicting human-supplied data
- ▶ **10→10,000 bits per sample**

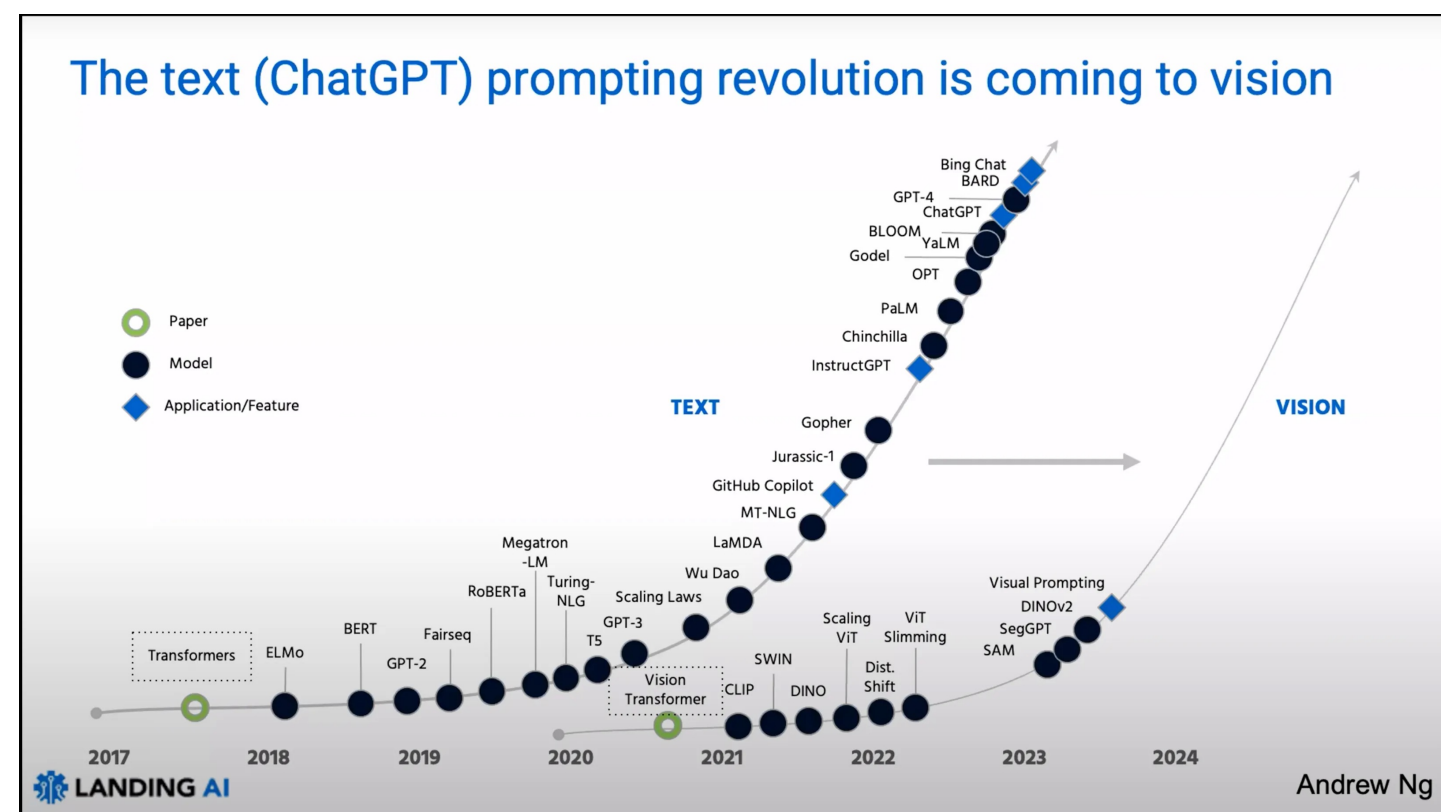
Unsupervised/Predictive Learning (cake)

- ▶ The machine predicts any part of its input for any observed part.
- ▶ Predicts future frames in videos
- ▶ **Millions of bits per sample**



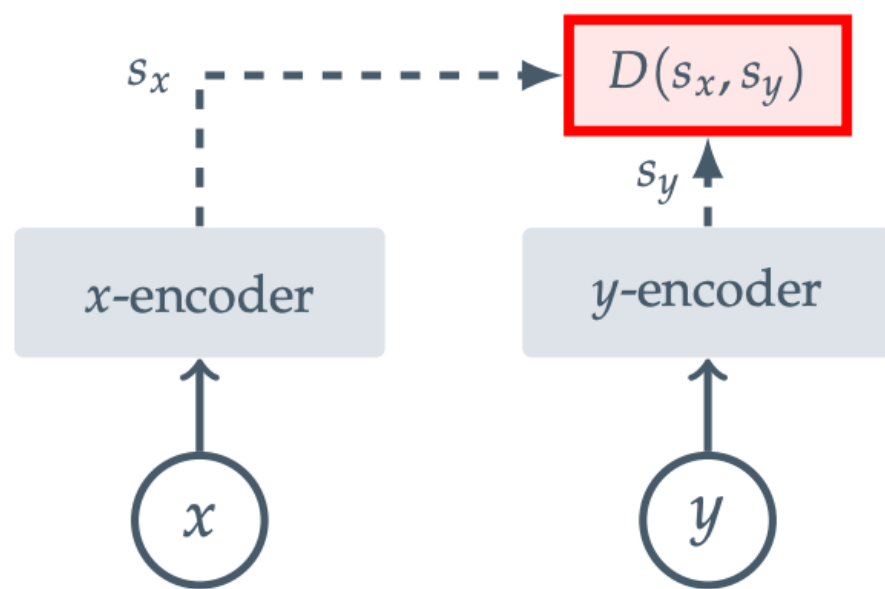
■ (Yes, I know, this picture is slightly offensive to RL folks. But I'll make it up)

Yann LeCun, ~2016

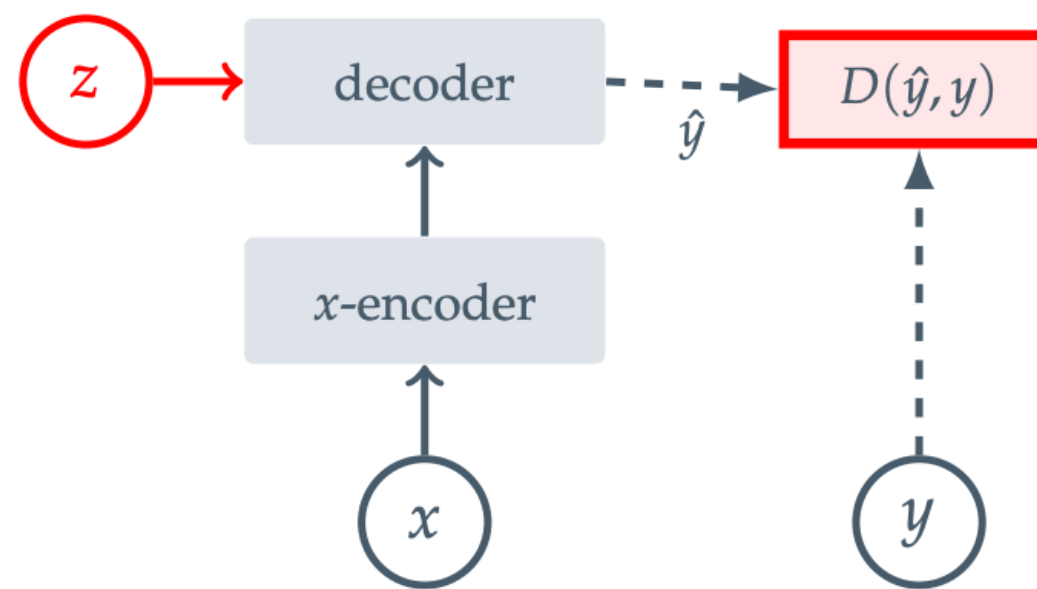


Andrew Ng, 2023

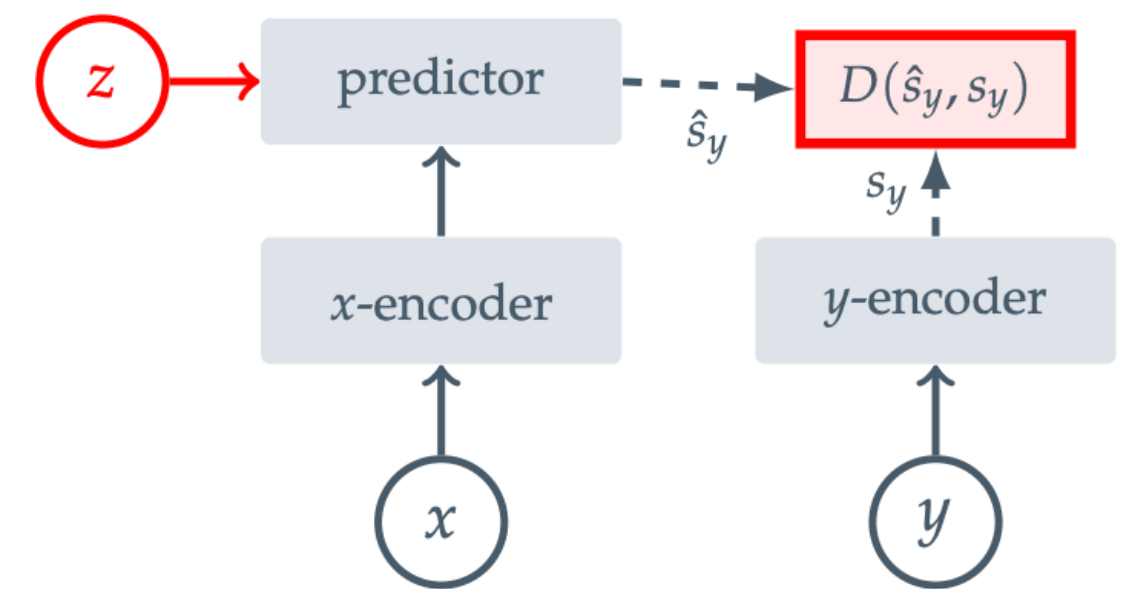
Core principle



(a) **Joint-Embedding Architecture**



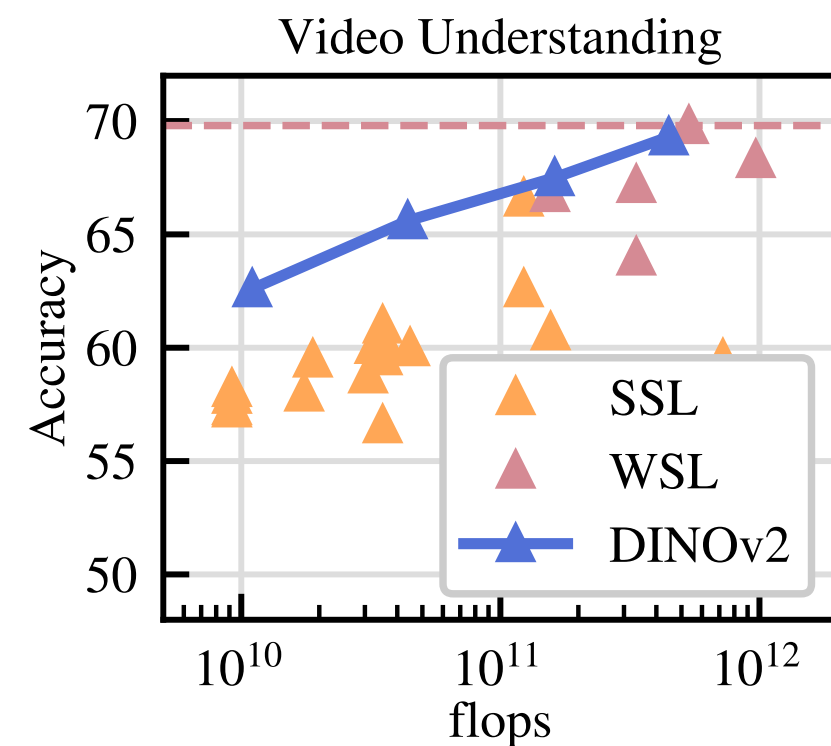
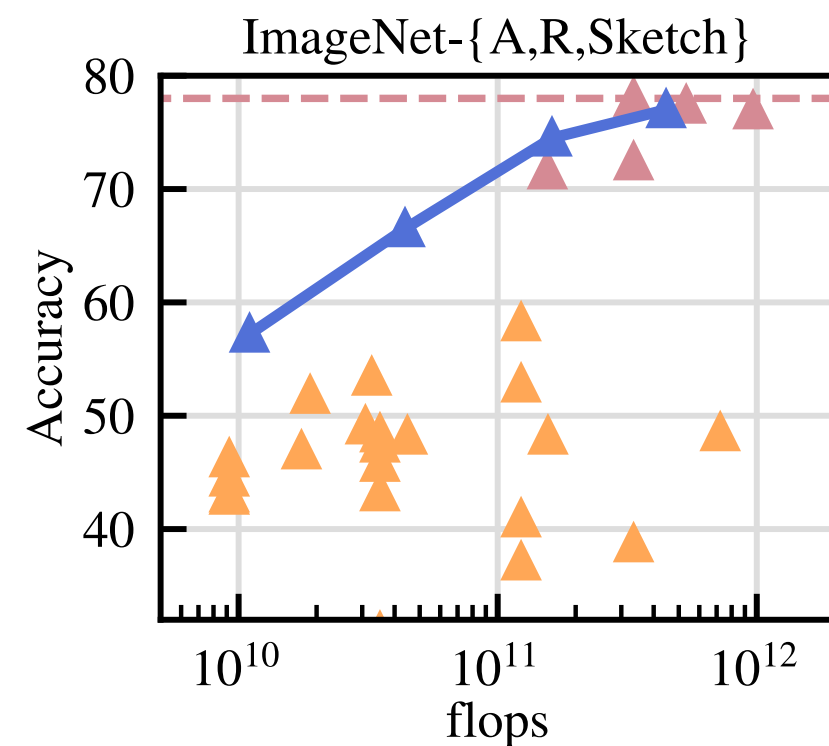
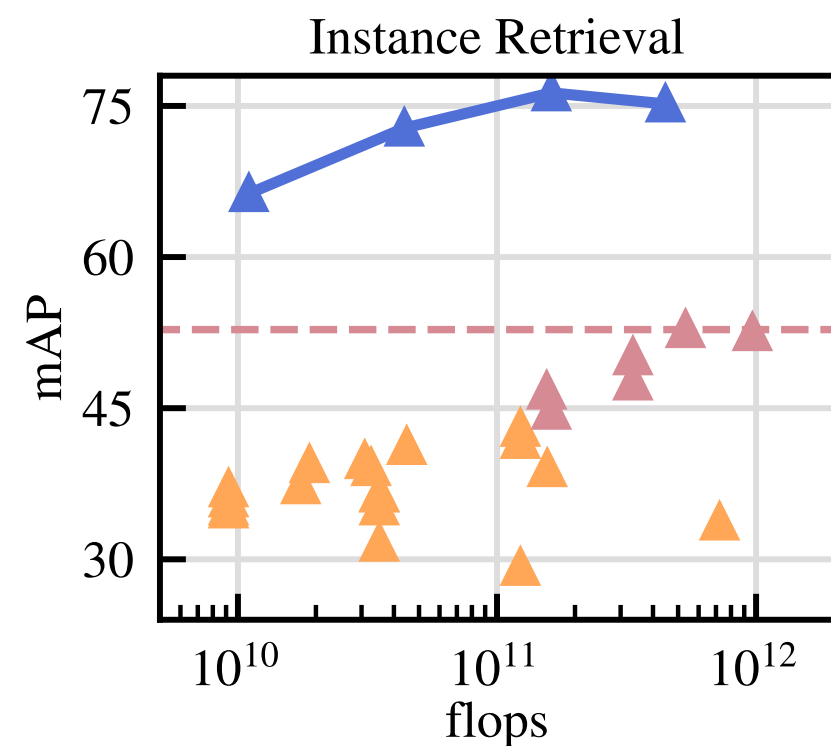
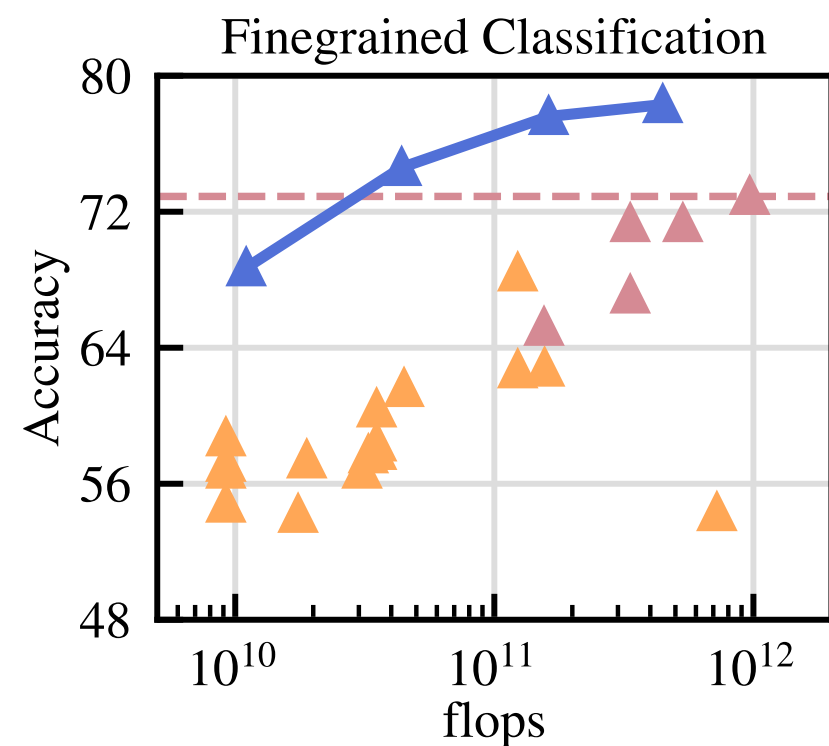
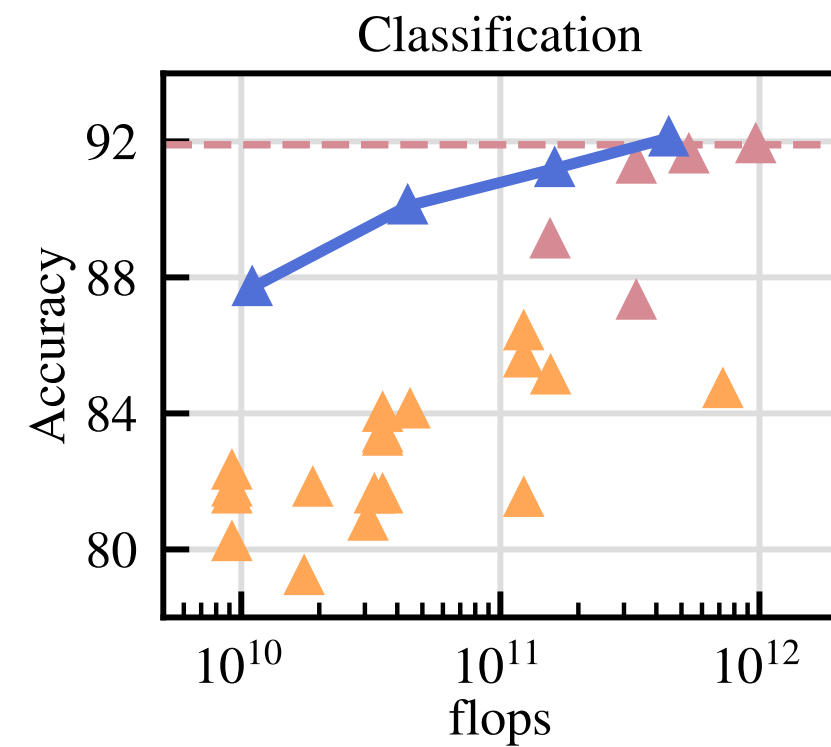
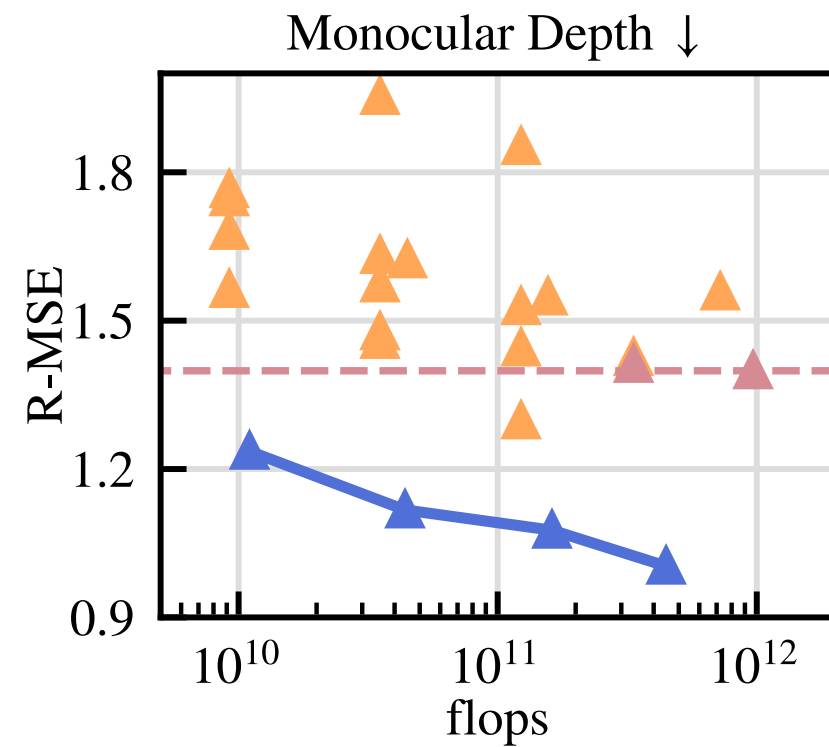
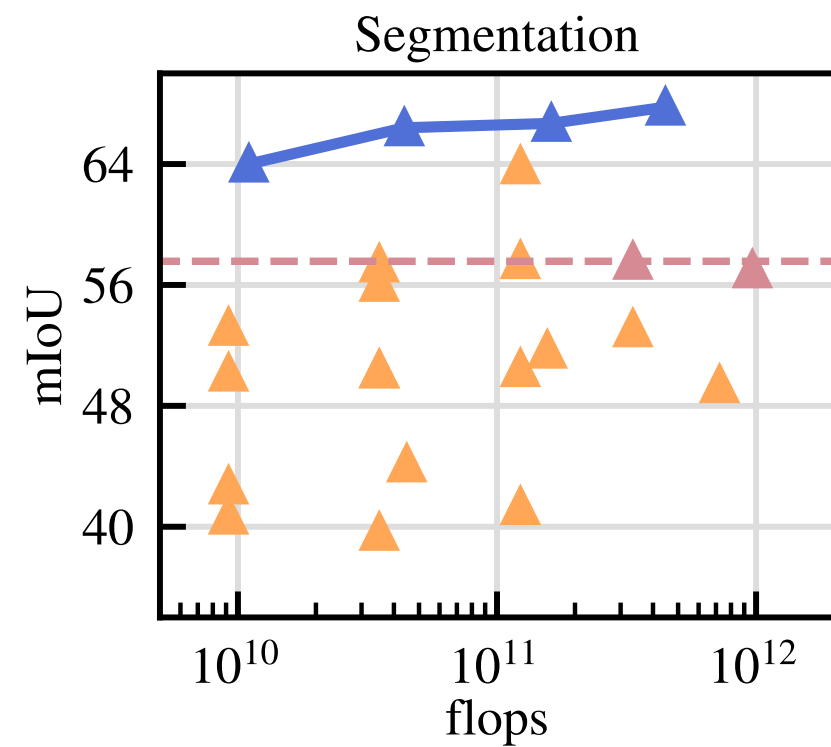
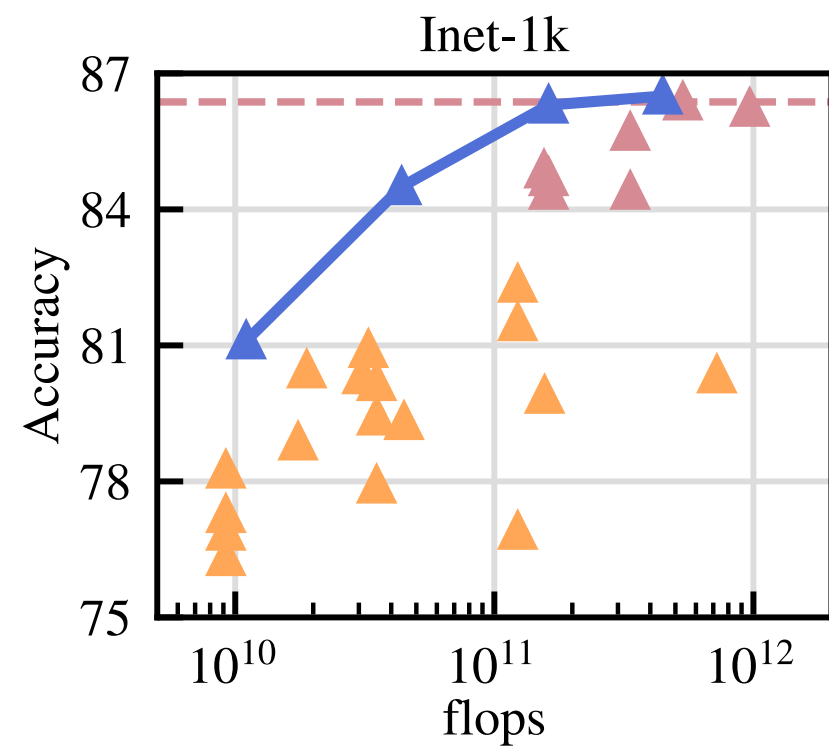
(b) **Generative Architecture**



(c) **Joint-Embedding Predictive Architecture**

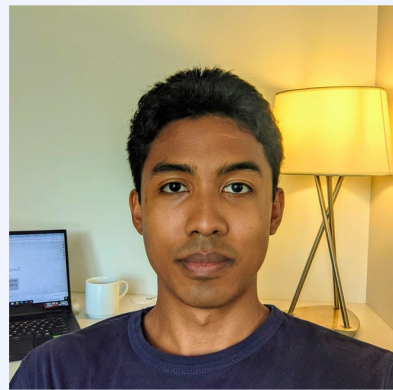
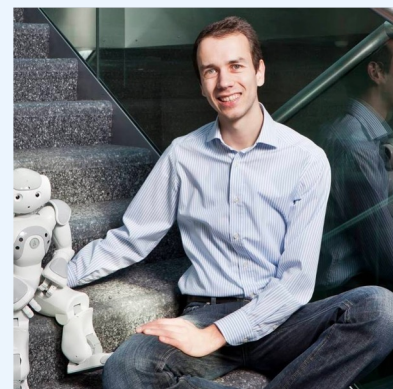
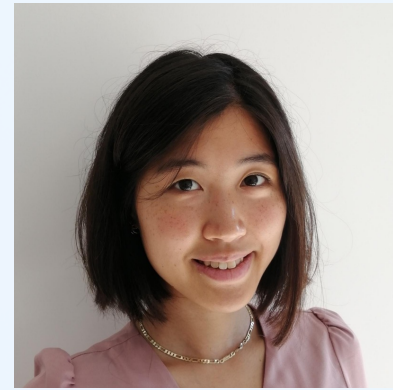
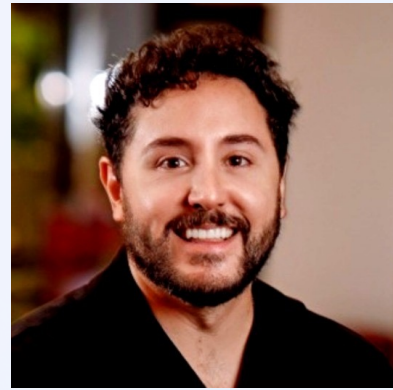
SplitBrainAE, NAT,
CPC, DeepCluster,
CPCv2, SELA, MoCo,
PIRL, SimCLR,
SwAV, MoCov2,
PCL, BYOL,
Barlow Twins, DINO,
SimCLRv2,
NN-CLR, VicReg,
BEiT, MAE...





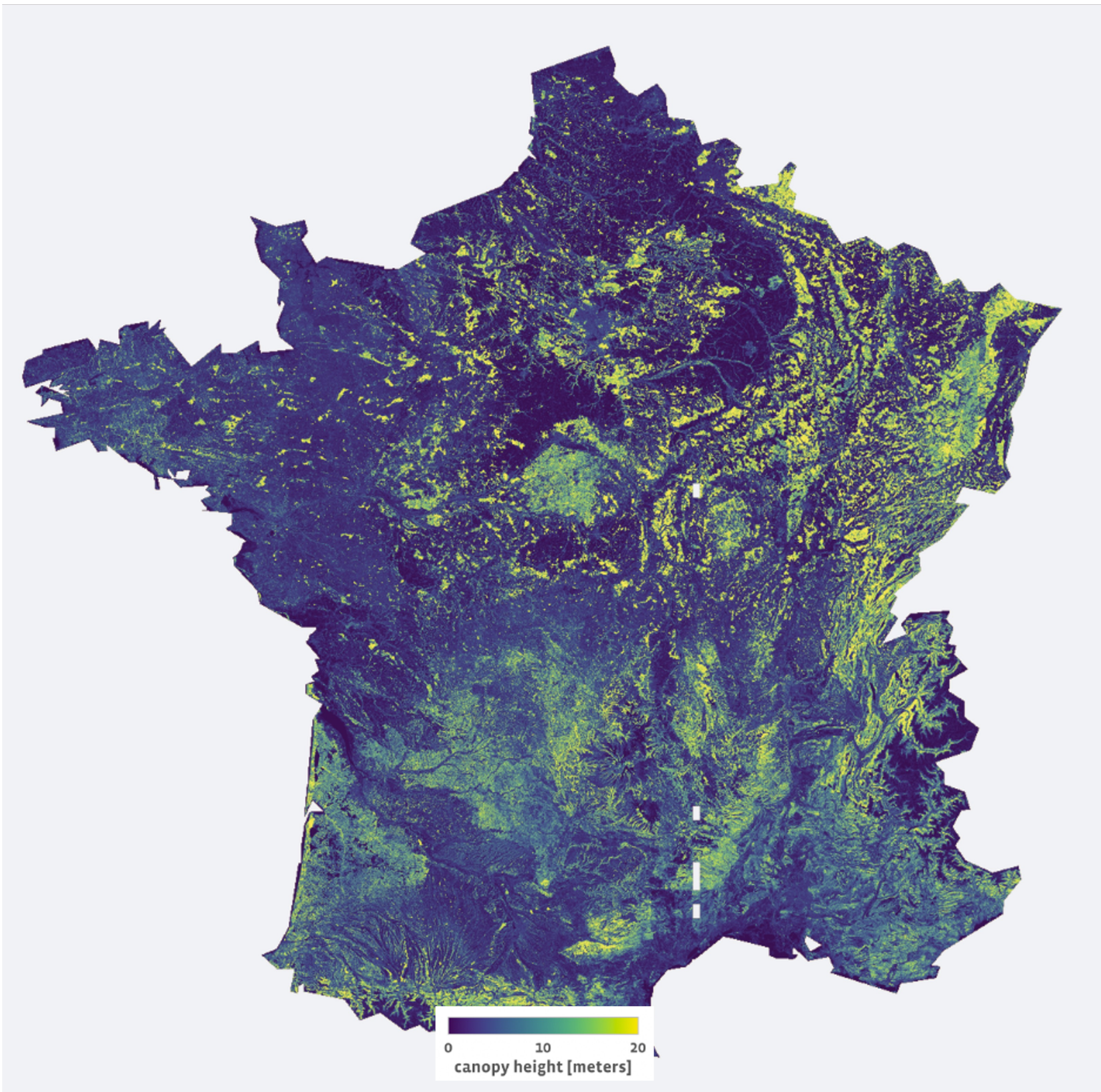
Scientific problems : a perfect application





Applications of SSL

High-Resolution Canopy Height Estimation



Physical Modelling @ Meta



Jamie Tolan



Ben Nosarzewski



Tobias Tiecke

World Ressource Institute



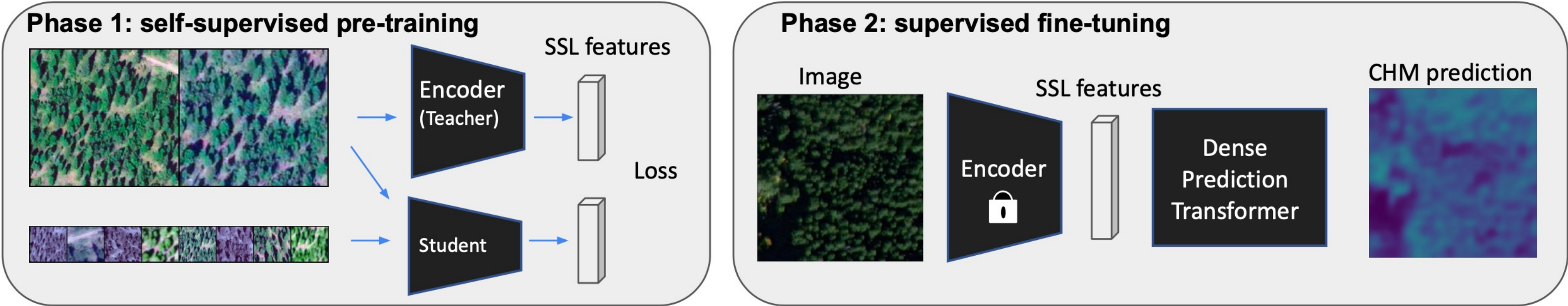
John Brandt



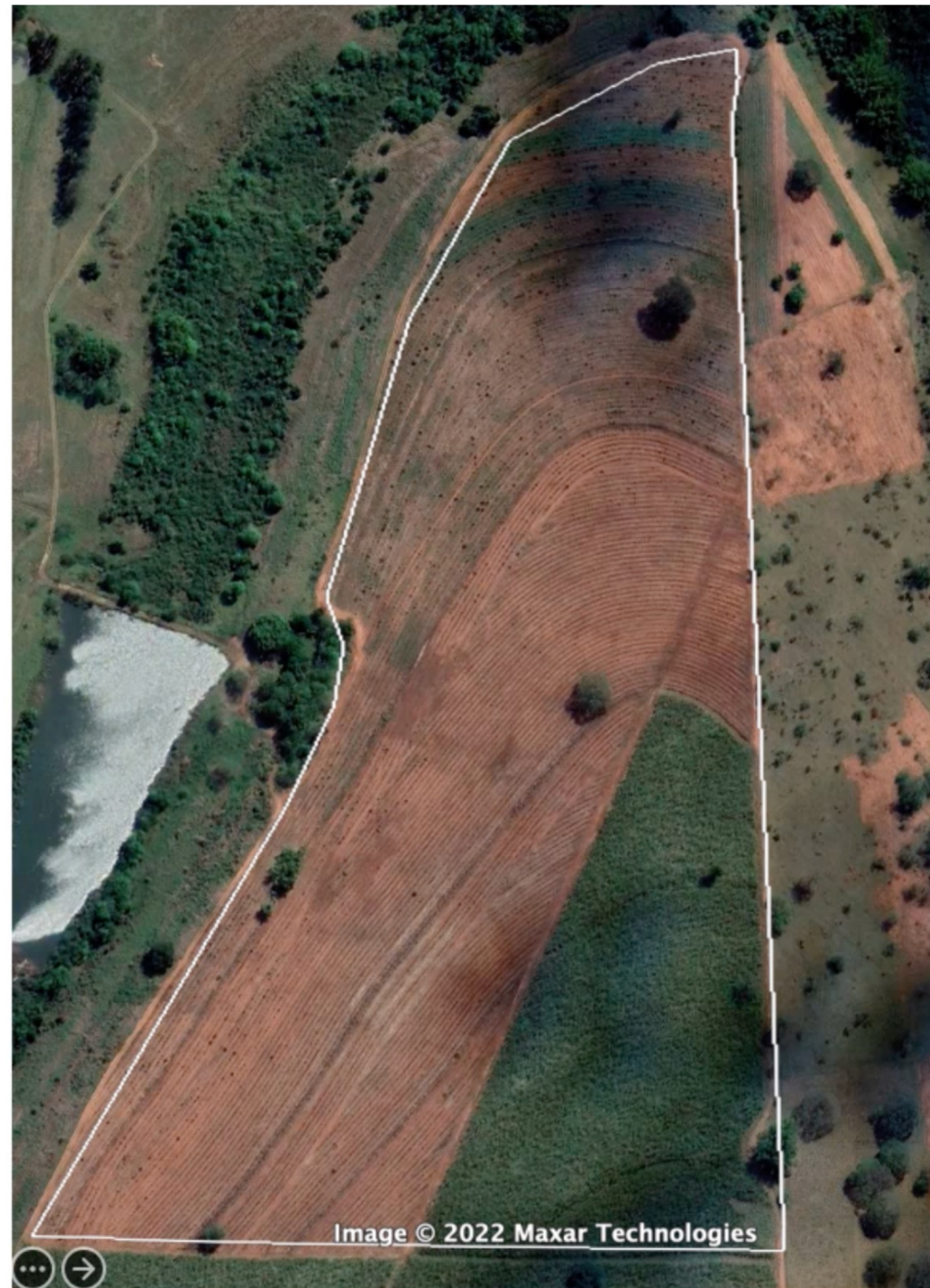
Justine Spore

Canopy Height Estimation

	Coverage	Type	Channels	Beam
MAXAR	Global	Satellite	RGB	0.5 m
GEDI	Near-Global	Satellite	RGB + LIDAR	25 m
NEON	Small	Airborne	RGB + LIDAR	1m



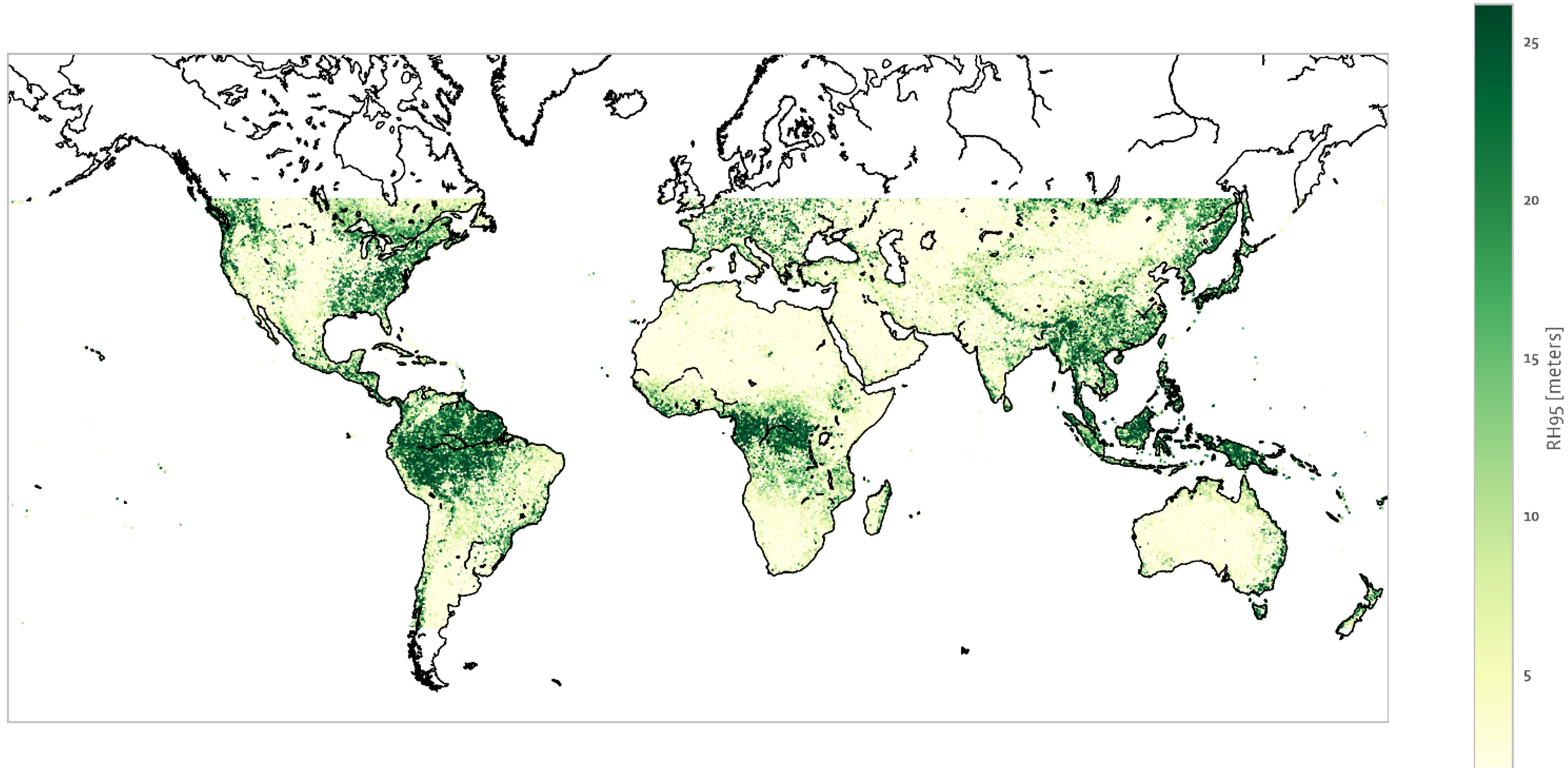
Worldwide Satellite orthorectified map



- 0.5 meter RGB stitched images
- Changes in terrain slope, view angle, sun angle, season, etc.

Datasets (Ground truth 1)

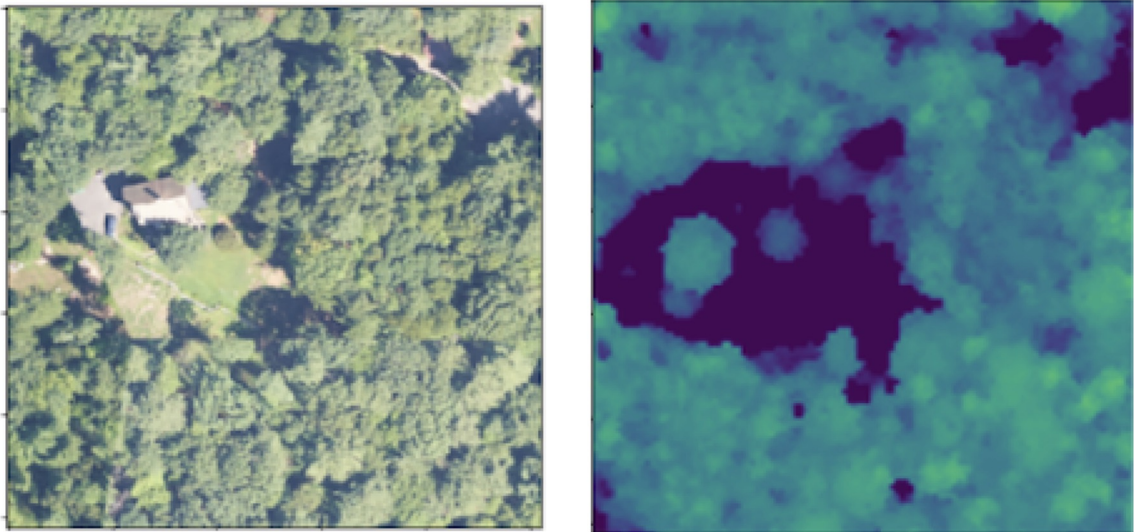
[GEDI](#): Space based lidar: Global (4% area coverage), 25 meter beamwidth, 1064nm, 1.28B data points:



Neon dataset (Ground truth 2)

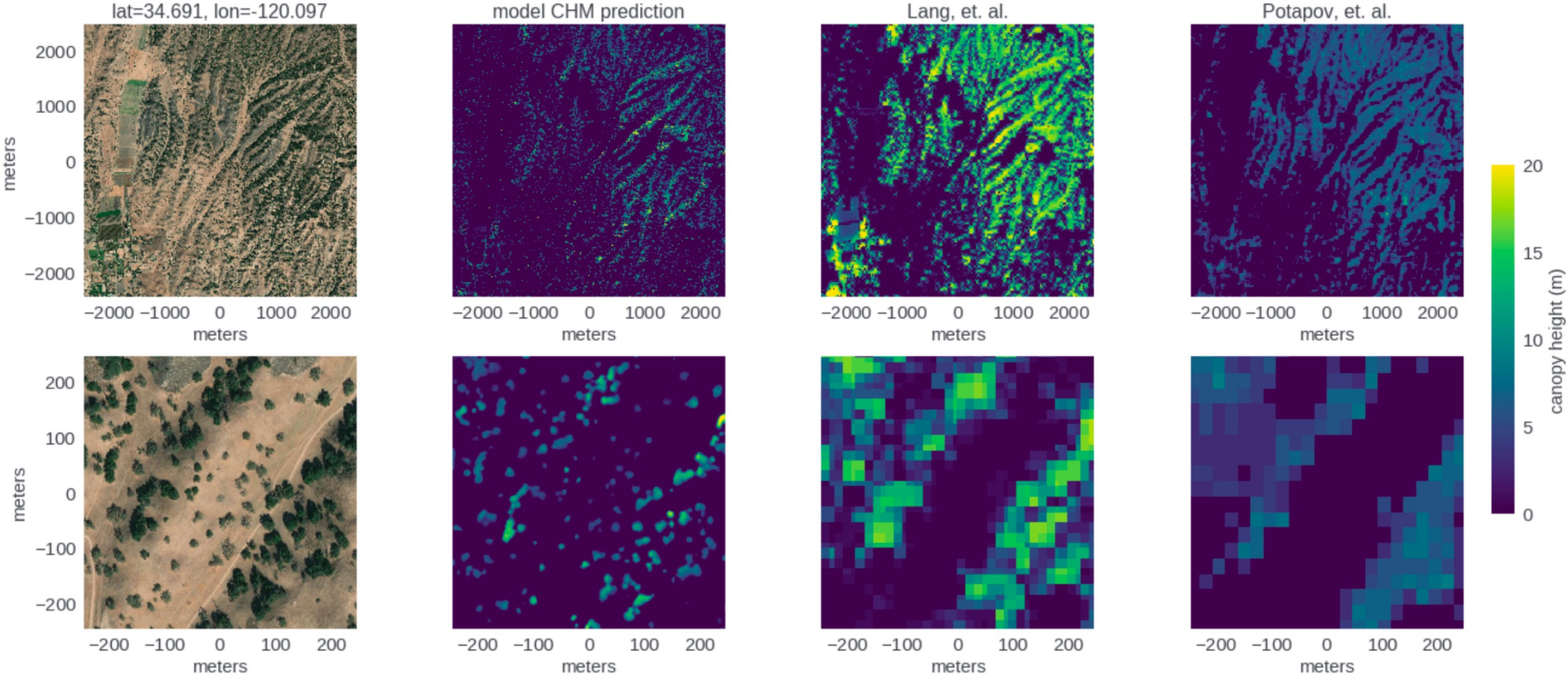
- Aerial Lidar canopy height maps
- 38 sites across North America
- About 5000 training image / lidar pairs of resolution 2200 pixels
- Dataset setup:
 - 80%-10%-10% train/val/test split
 - 5 Test sites completely independent of training sites

2

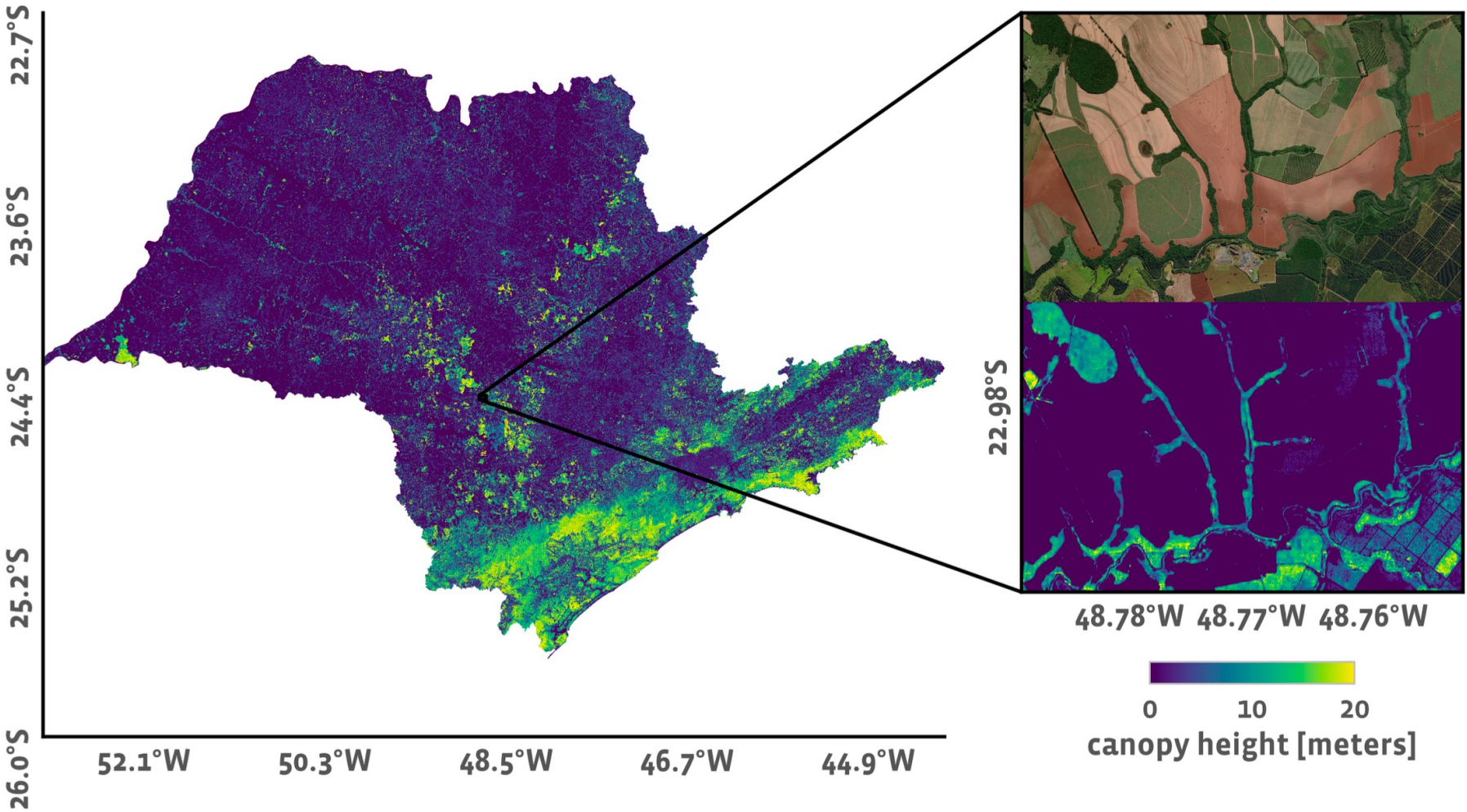
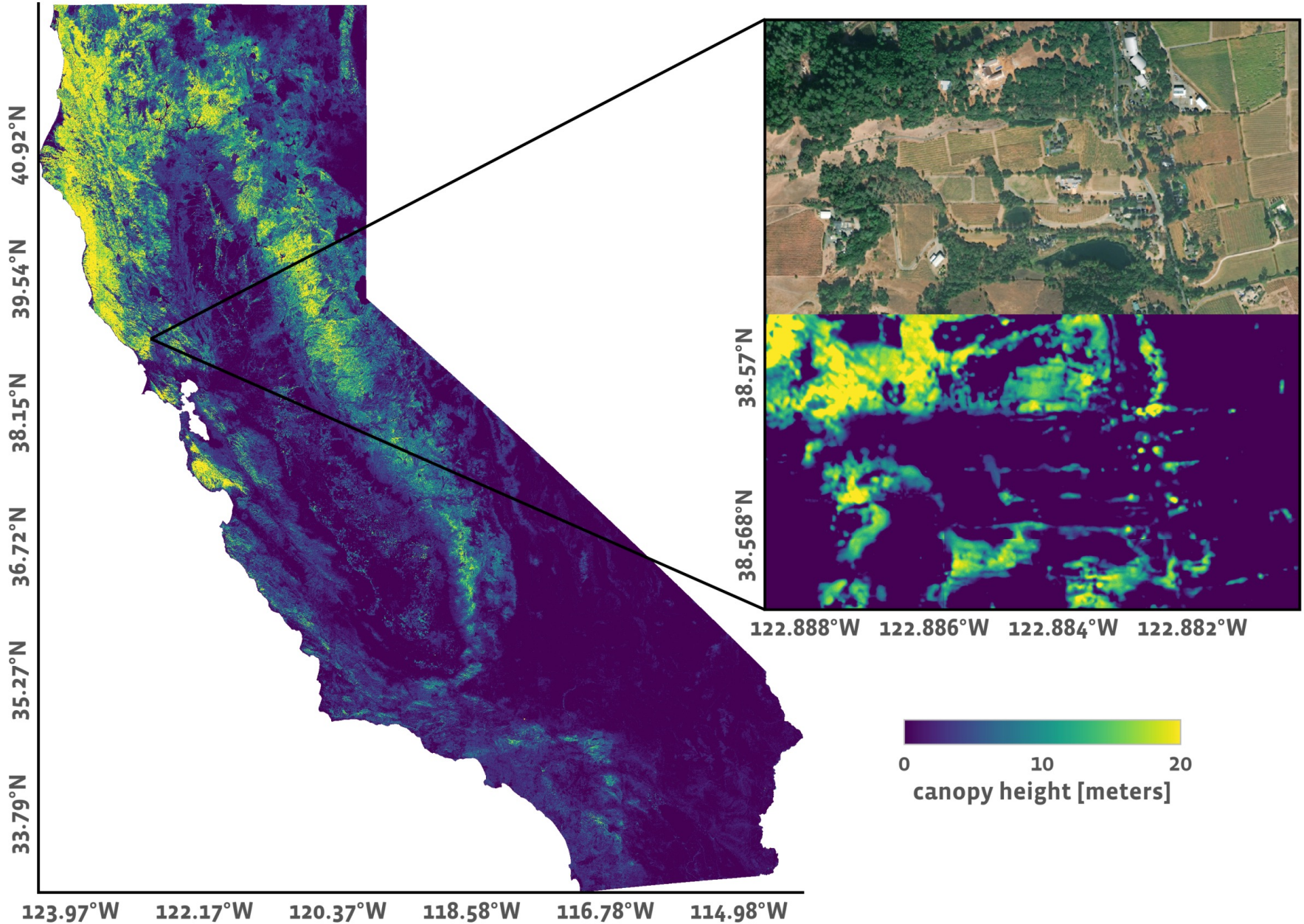


In addition to the Neon dataset, we used some lidar data in Sao Paulo, California and France for evaluation.

Canopy Height Estimation

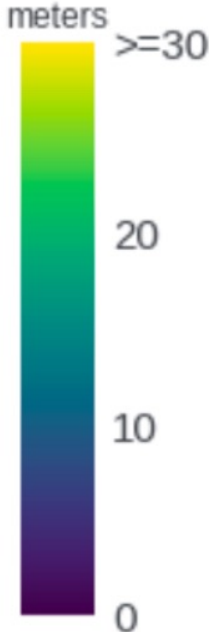
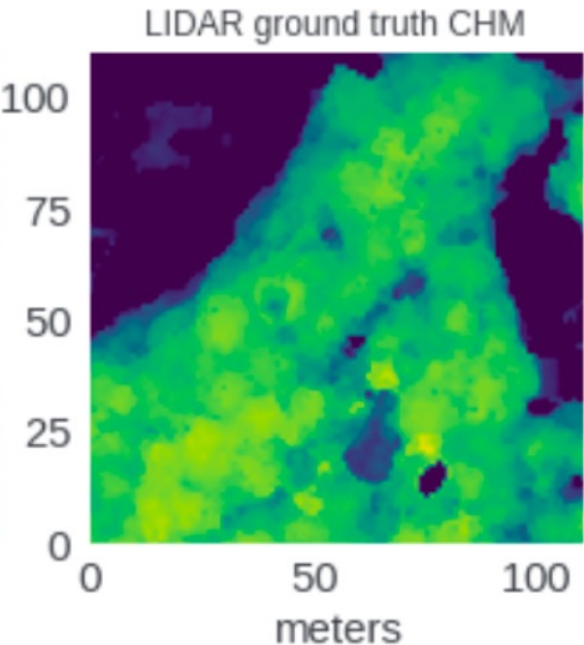
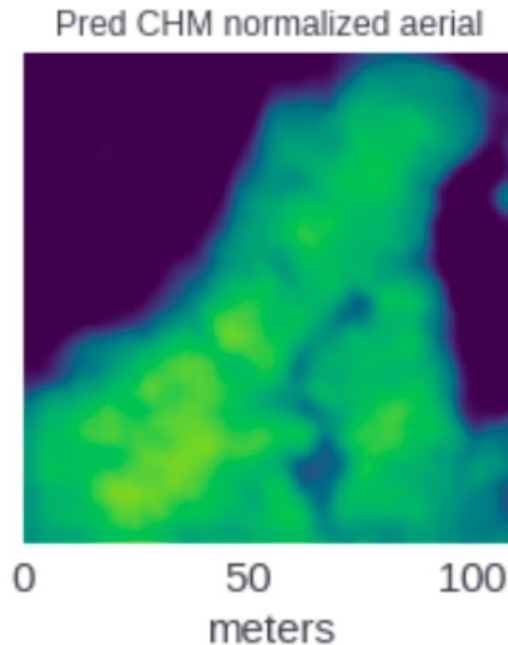
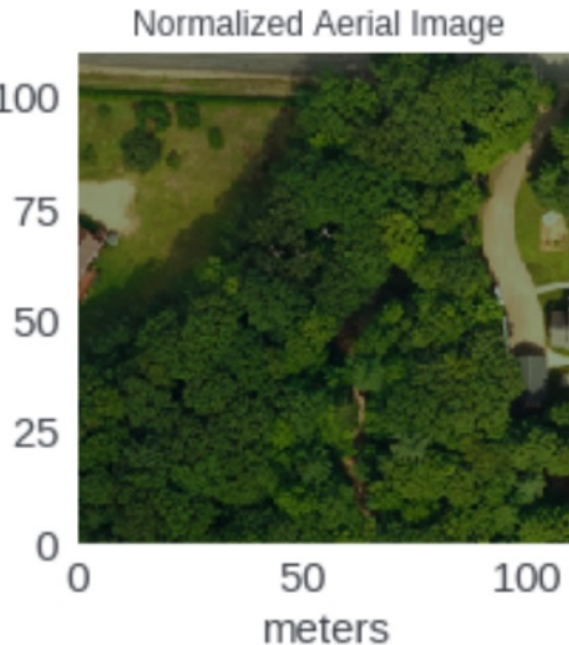
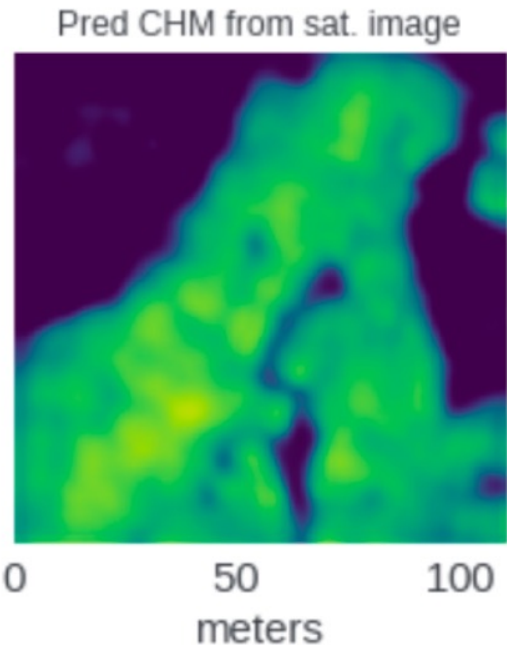
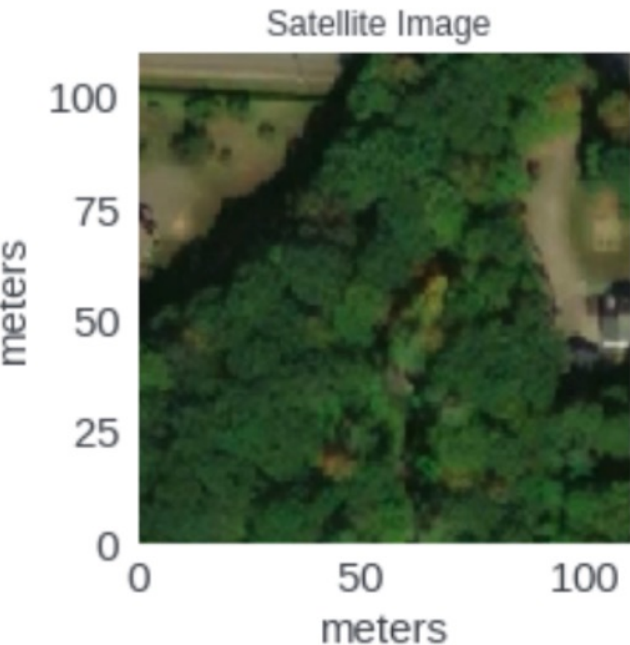


Canopy Height Estimation



Canopy Height Estimation

lat=44.083 lon=-71.296



Single-Cell Microscopy



Juan C. Caicedo
University of Wisconsin-Madison /
Broad Institute of MIT

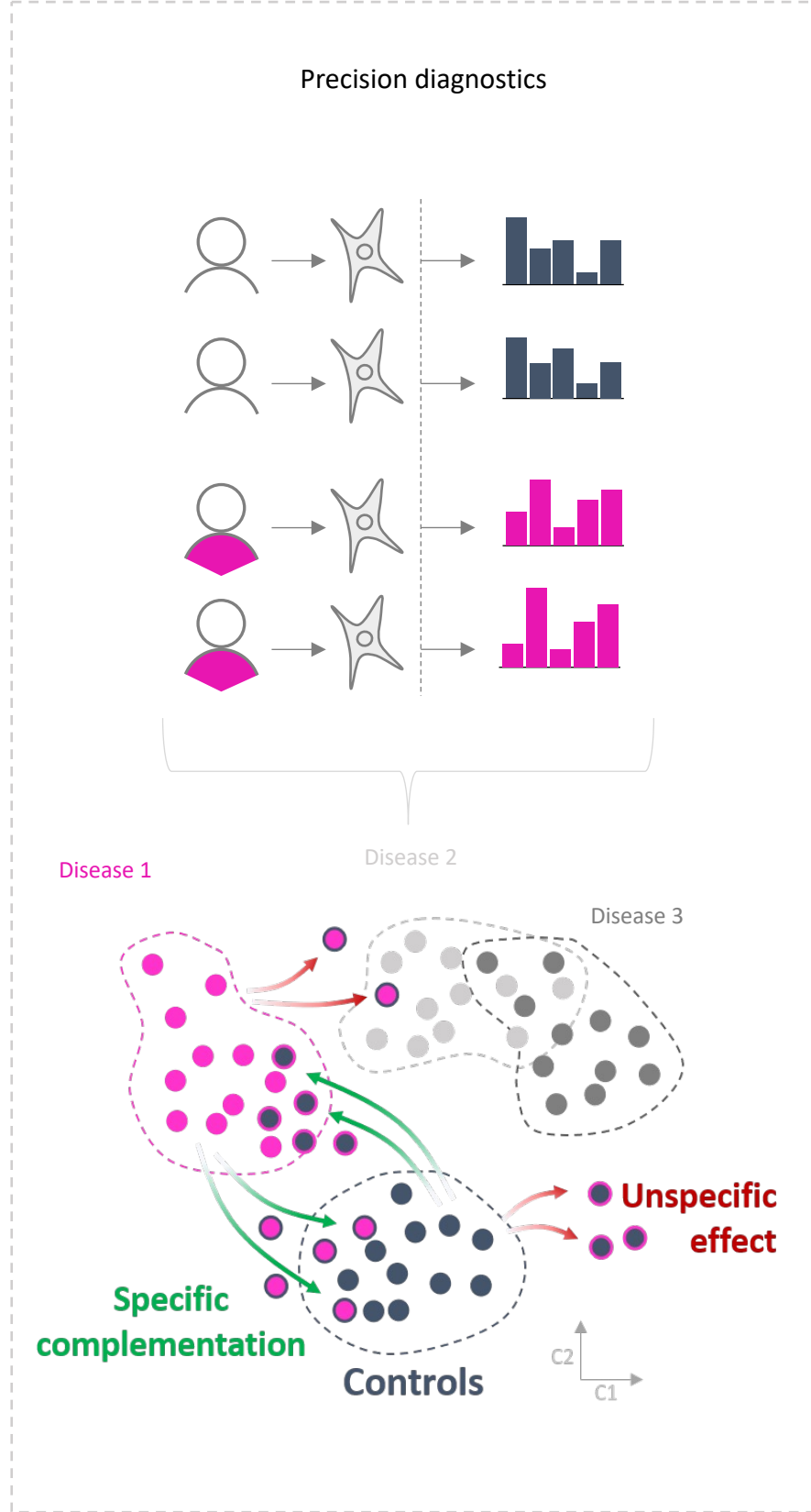
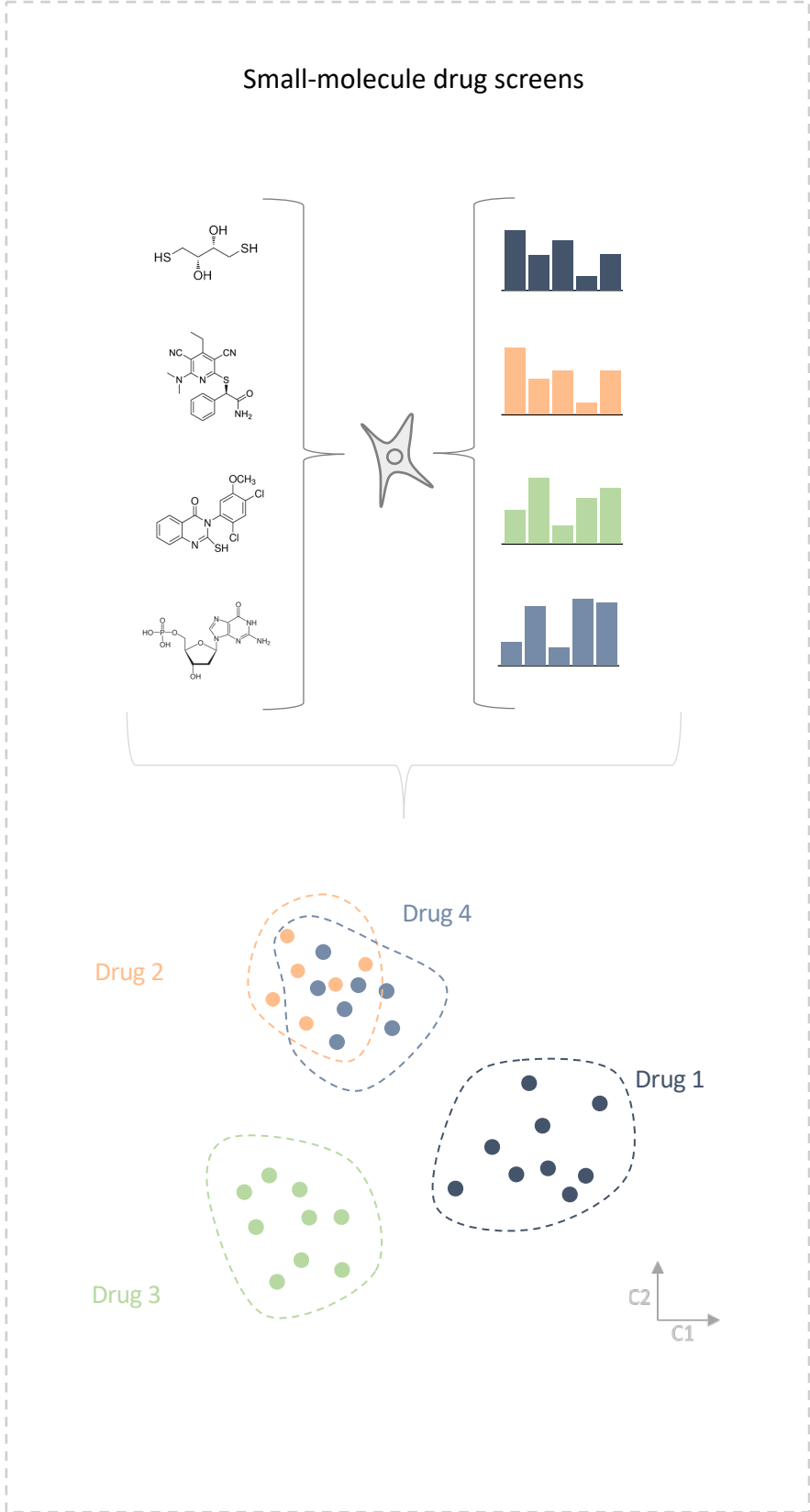
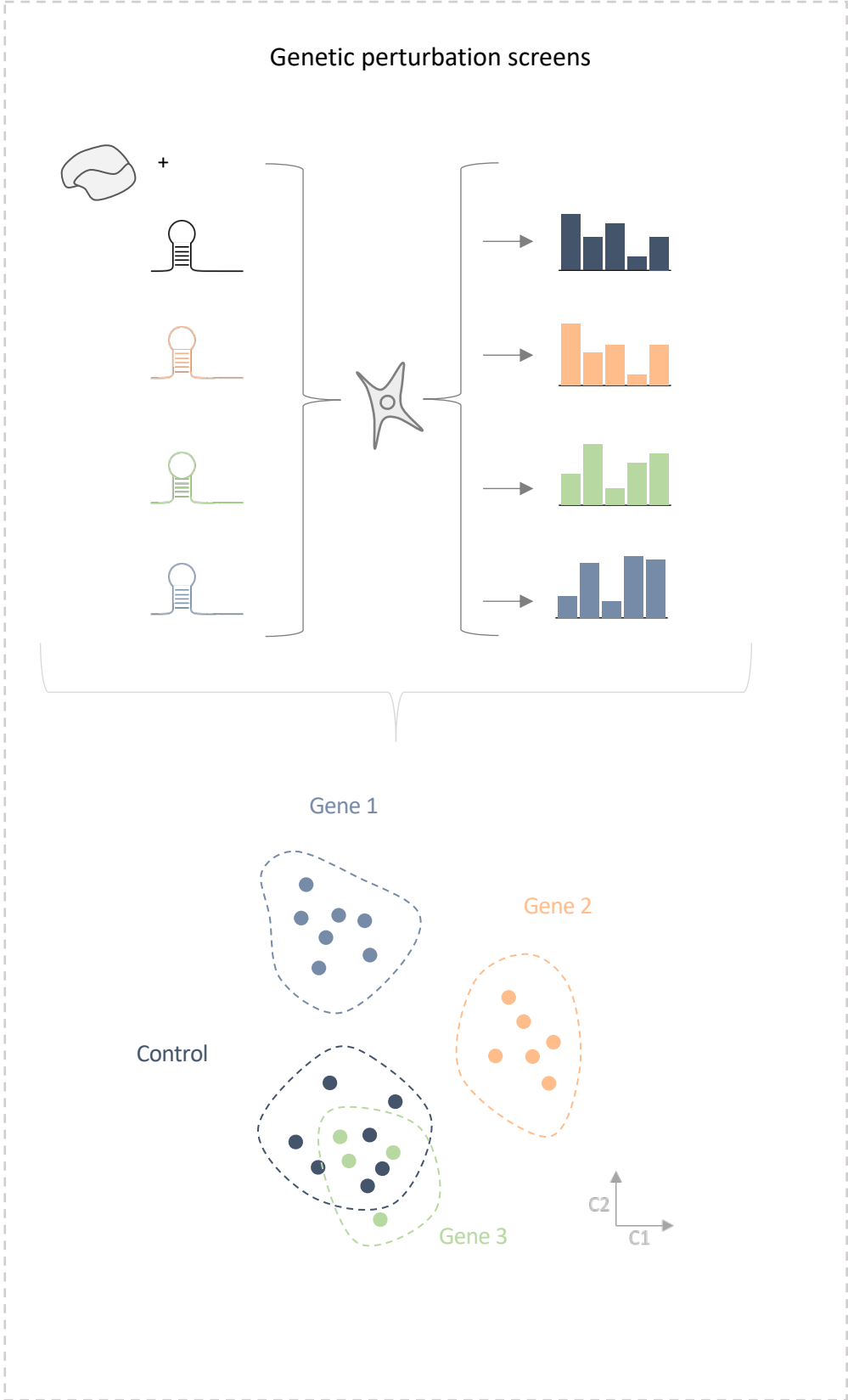


Wolfgang Pernice
Columbia University Irving Medical Center

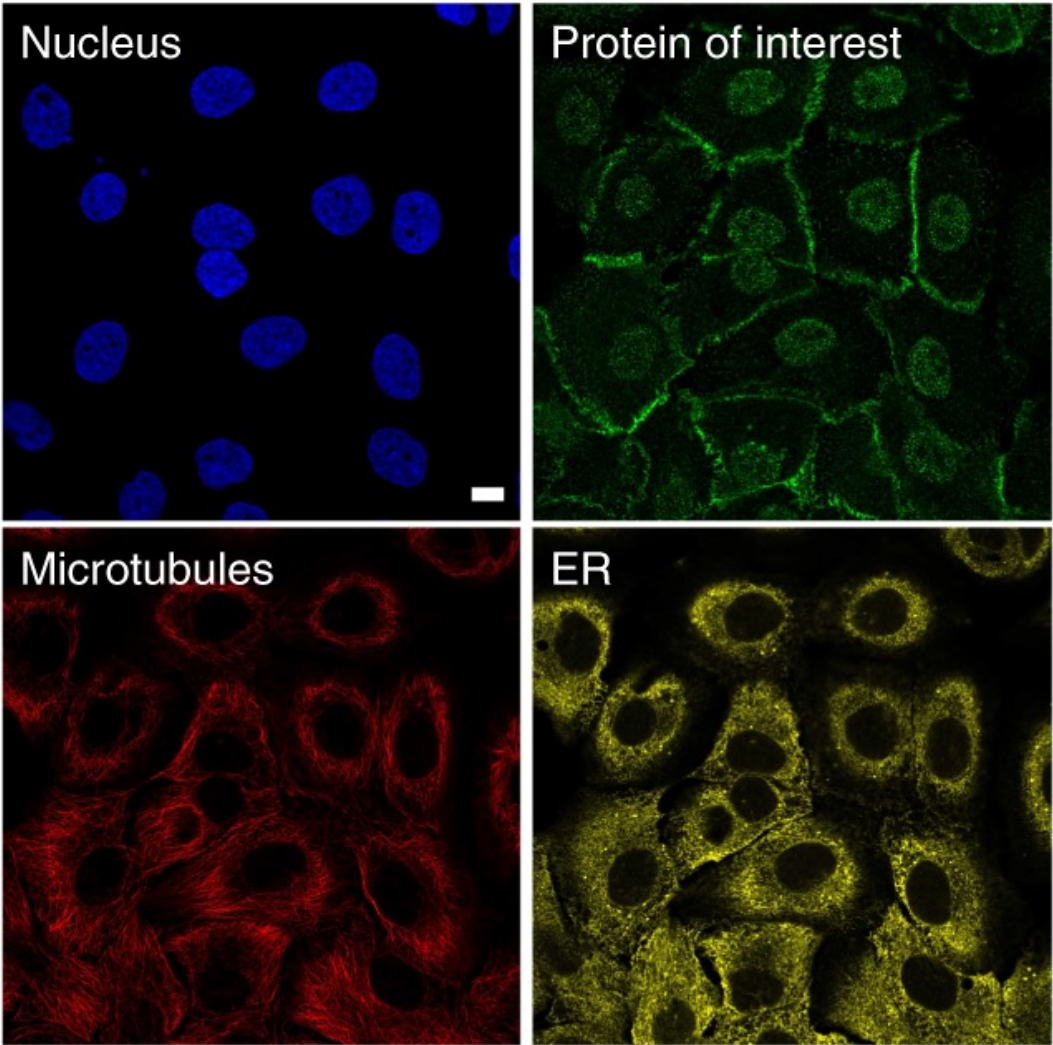


Michael Doron
Q.AI / Broad Institute of MIT

Single-Cell Microscopy



a



Classifier
→

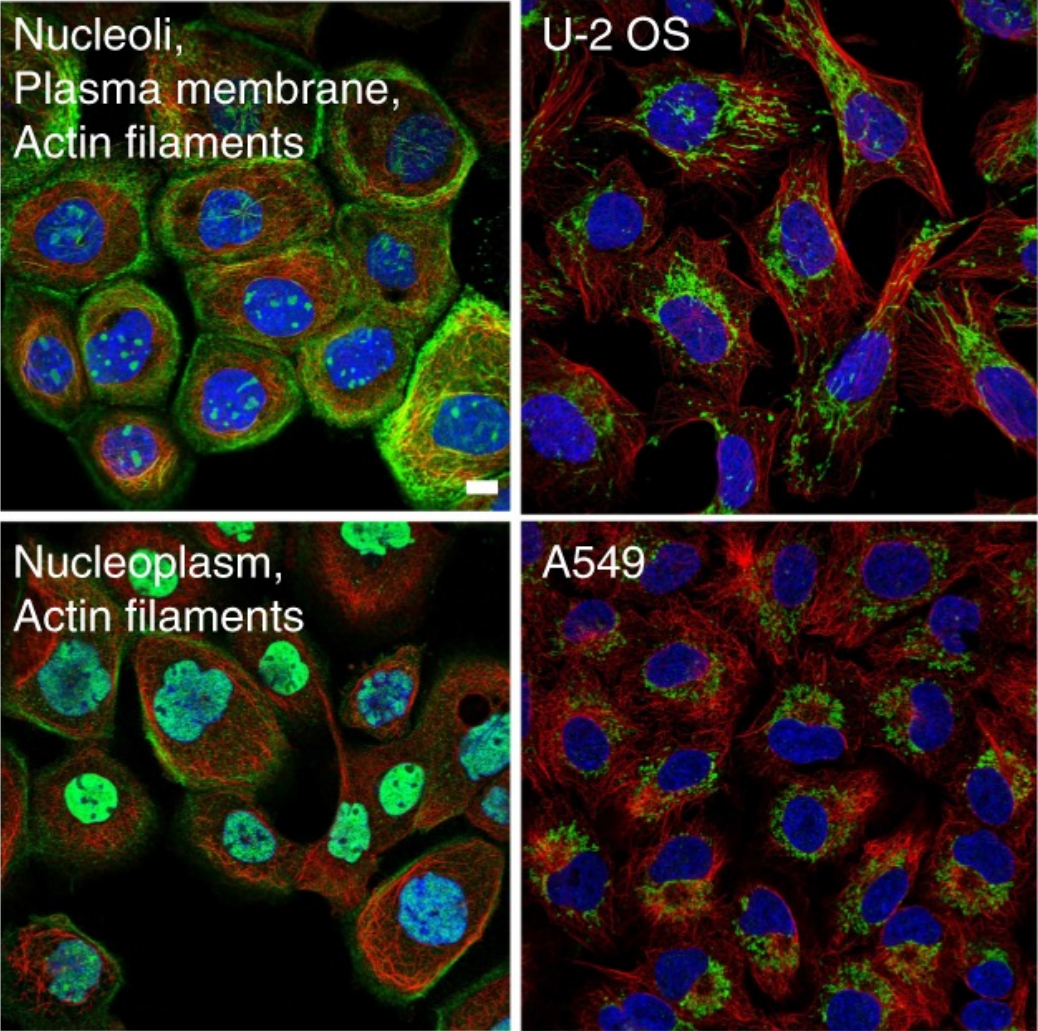
Multi-label prediction

Nucleoplasm
Cytosol
Plasma membrane
Nucleoli
Mitochondria
Golgi apparatus
Nuclear bodies
Nuclear speckles
Nucleoli fibrillar c.
Centrosome
Cell junctions
Actin filaments
...

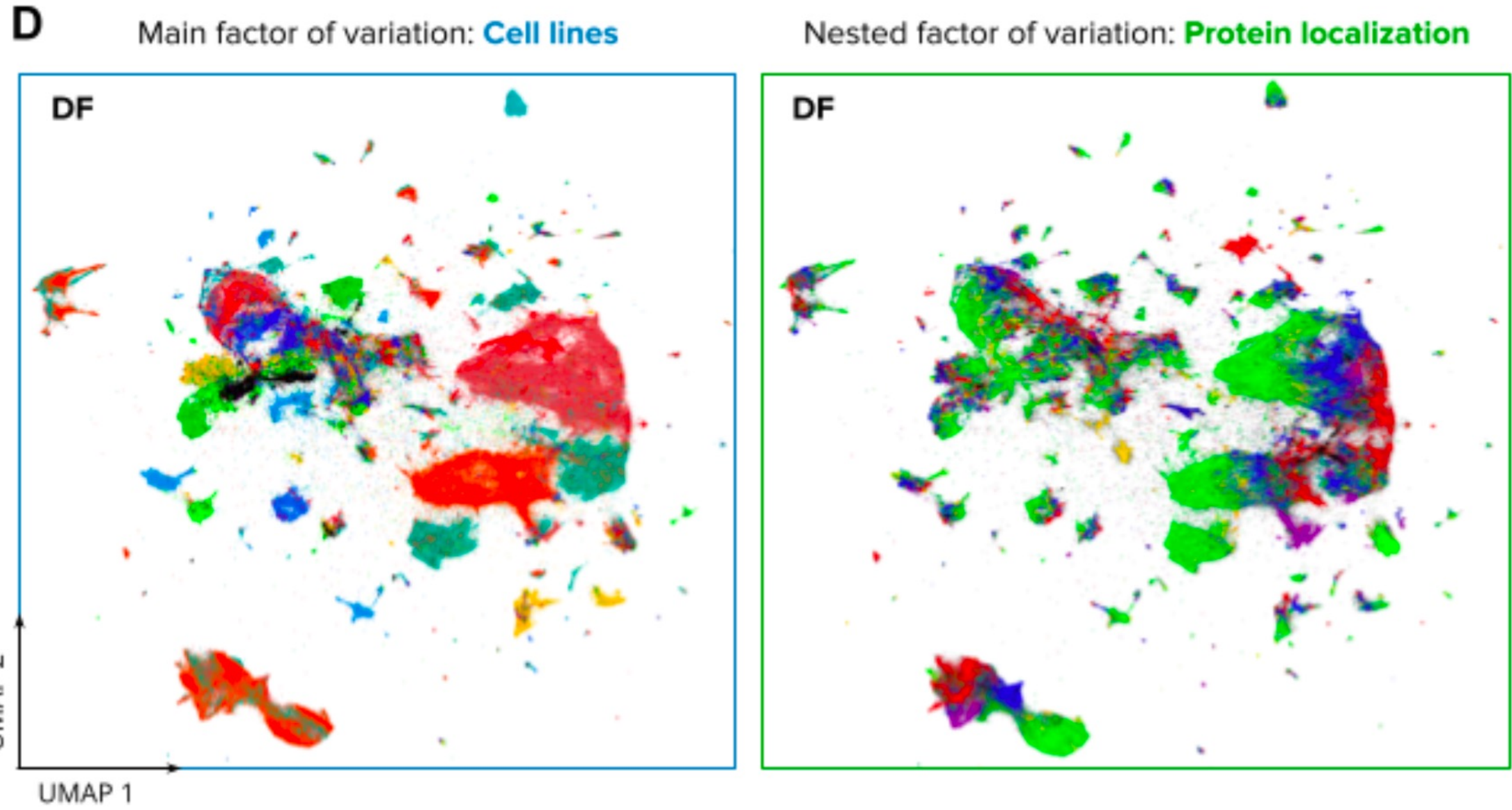
b

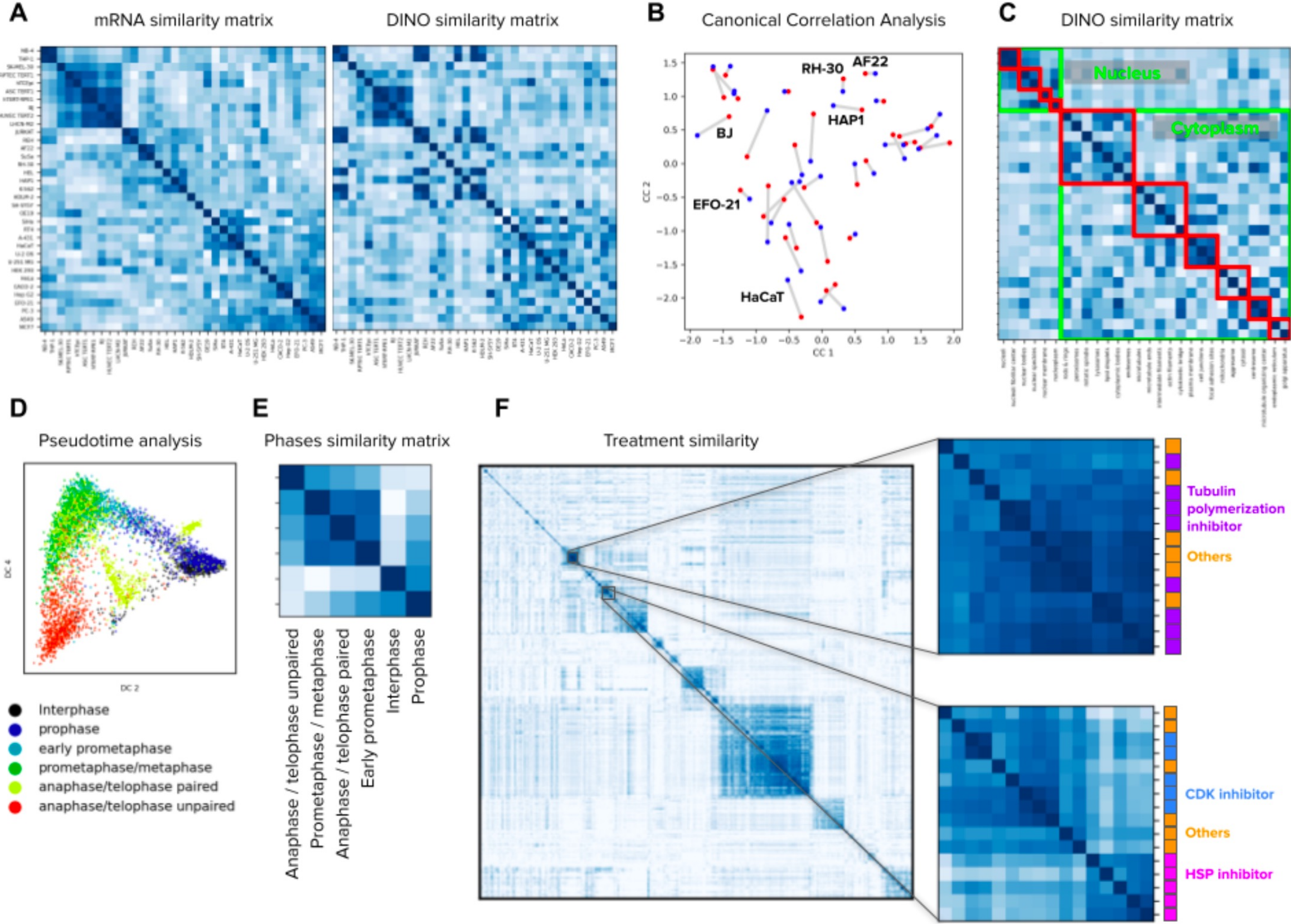
Multilocalizing proteins

Cell line variation



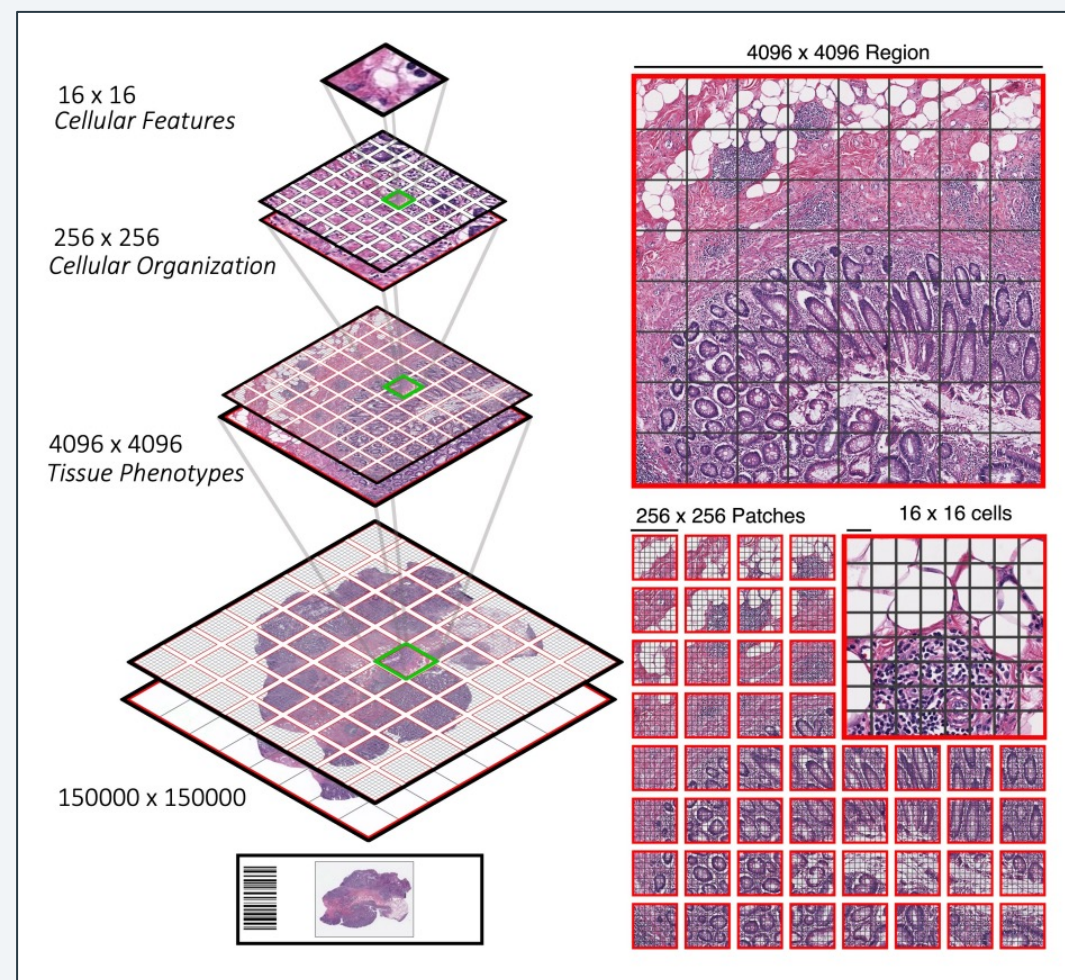
Uhlén, Mathias, et al. "Tissue-based map of the human proteome." *Science* 347.6220 (2015): 1260419.



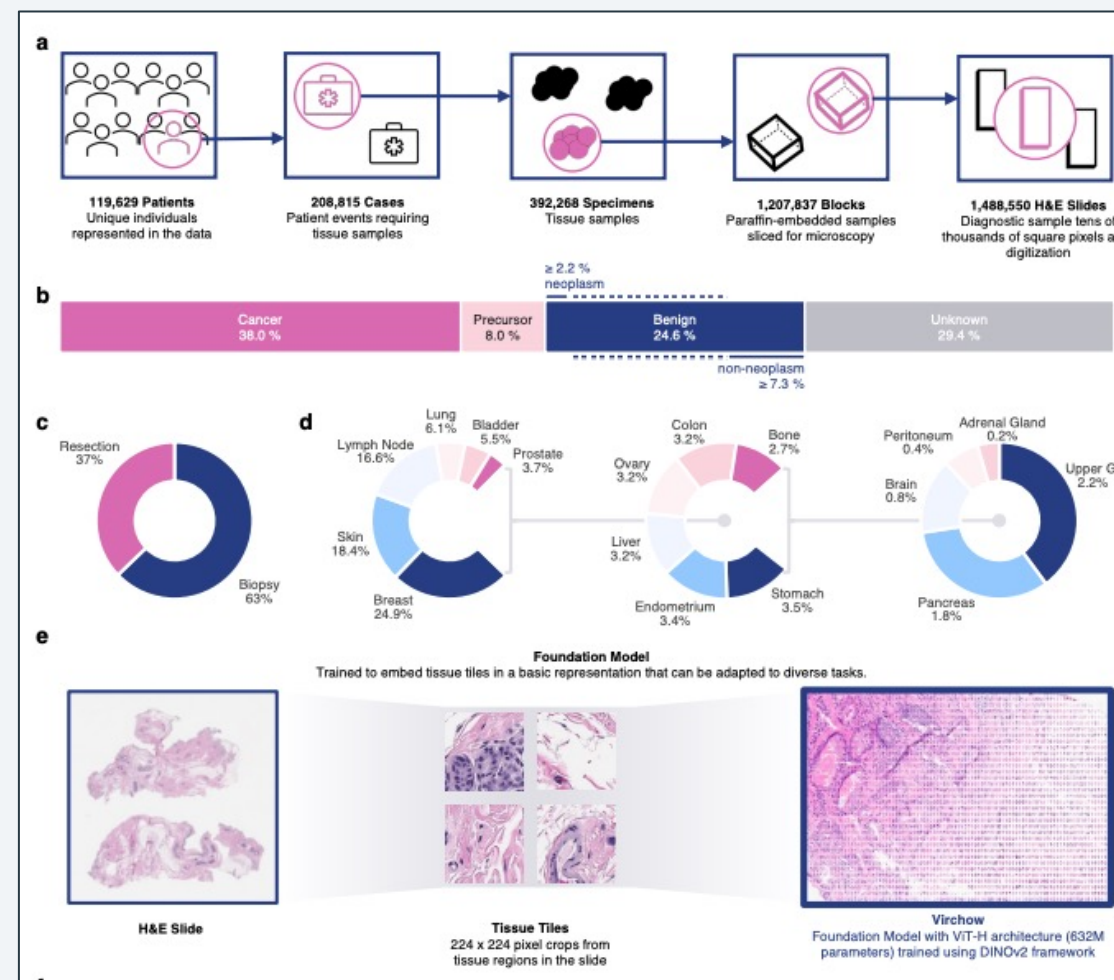


Conclusion and Future Work

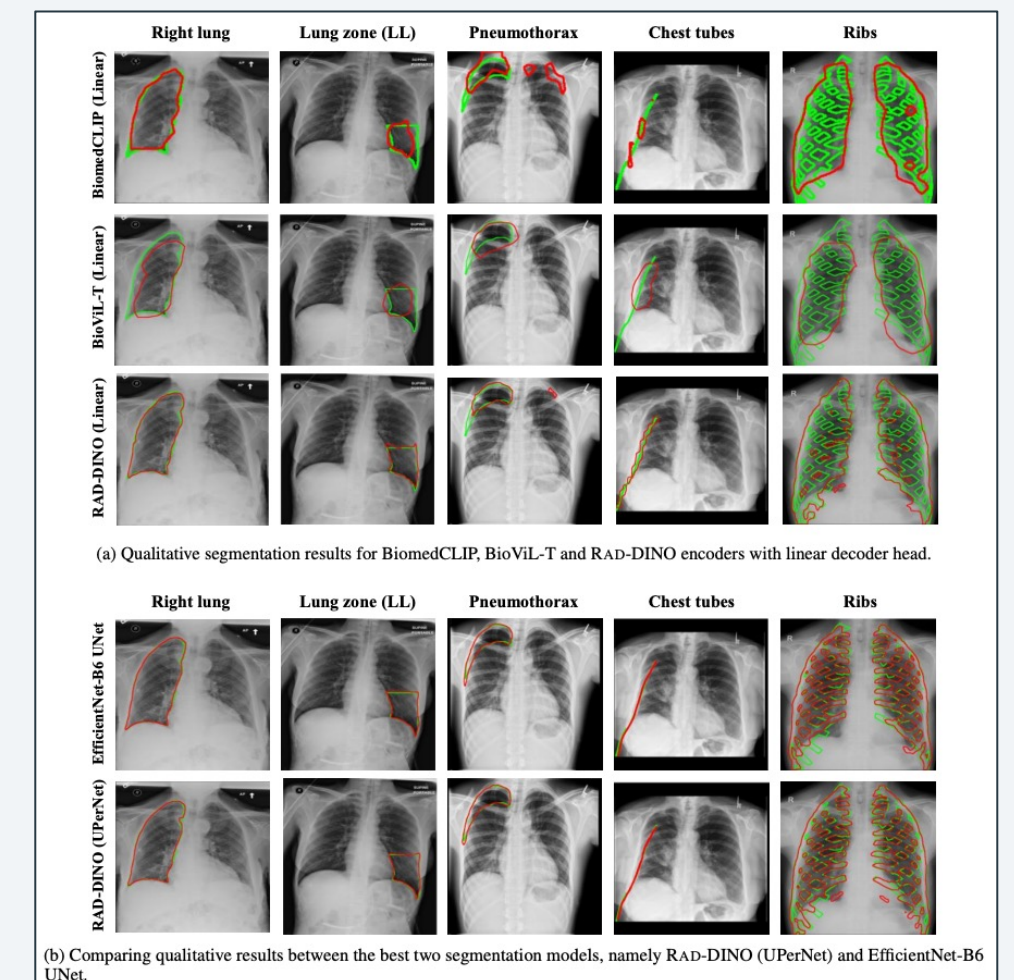
DINOv2 : a foundation model factory



Chen, Richard J., et al. "Scaling vision transformers to gigapixel images via hierarchical self-supervised learning." Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 2022.

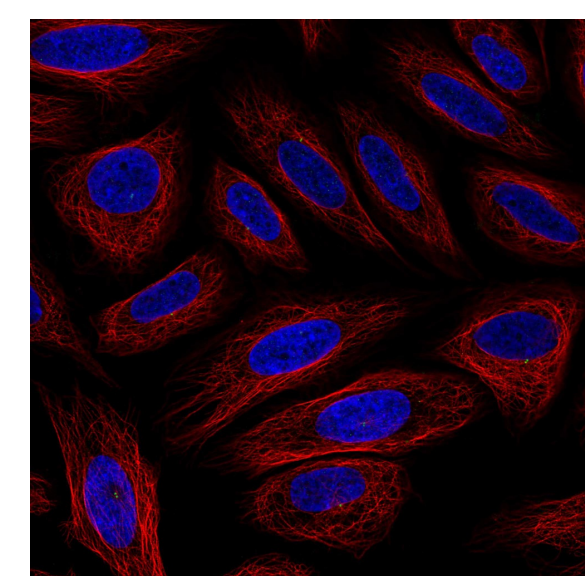
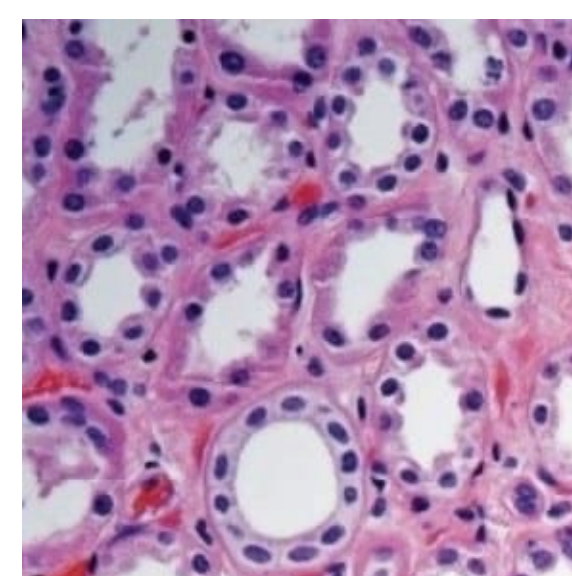
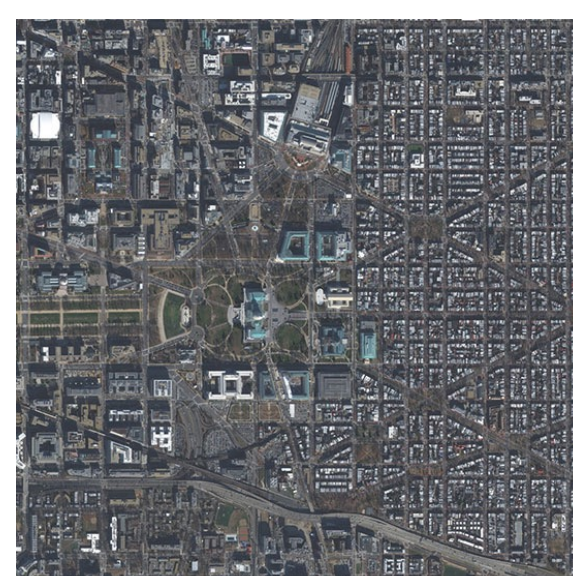
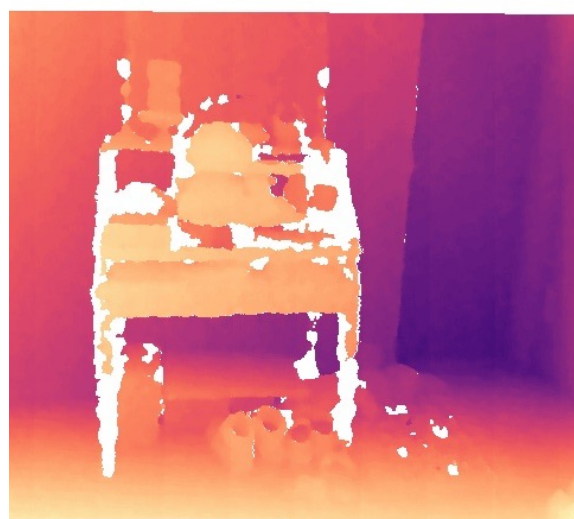
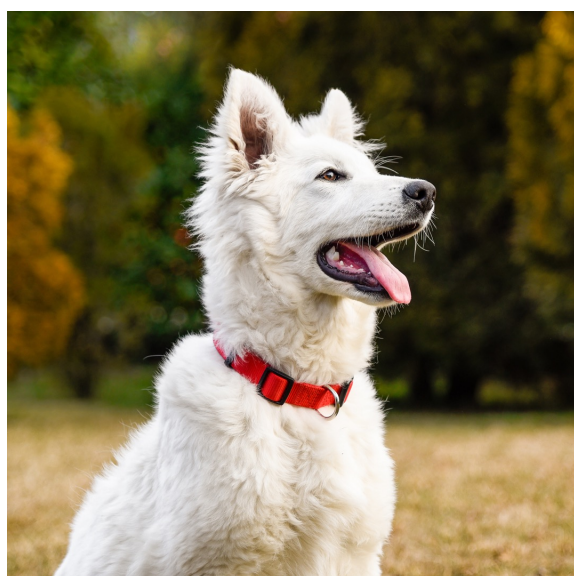
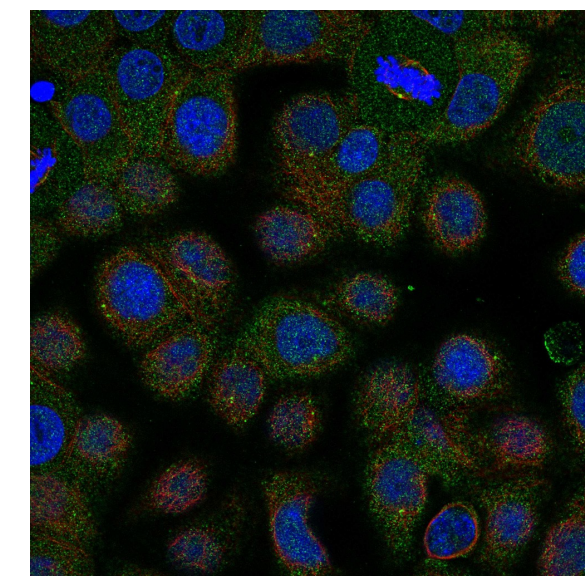
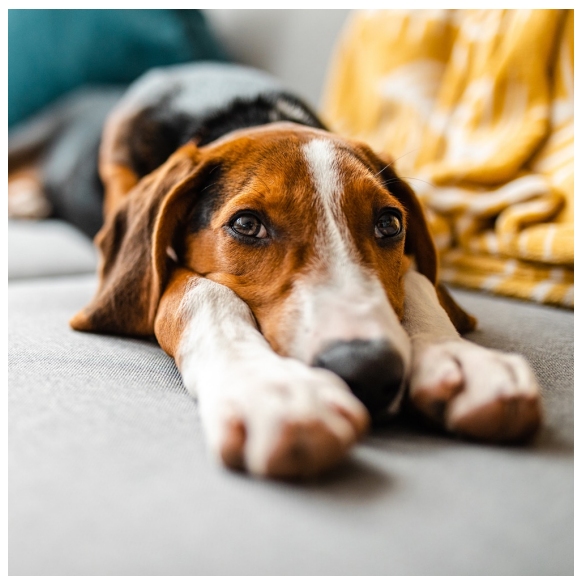


Vorontsov, Eugene, et al. "Virchow: A million-slide digital pathology foundation model." arXiv preprint arXiv:2309.07778 (2023).



Pérez-García, Fernando, et al. "RAD-DINO: Exploring Scalable Medical Image Encoders Beyond Text Supervision." arXiv preprint arXiv:2401.10815 (2024).

Learning Universal Visual Representations



 **Meta AI**