

ViTally Consistent: Scaling Biological Representation Learning for Cell Microscopy

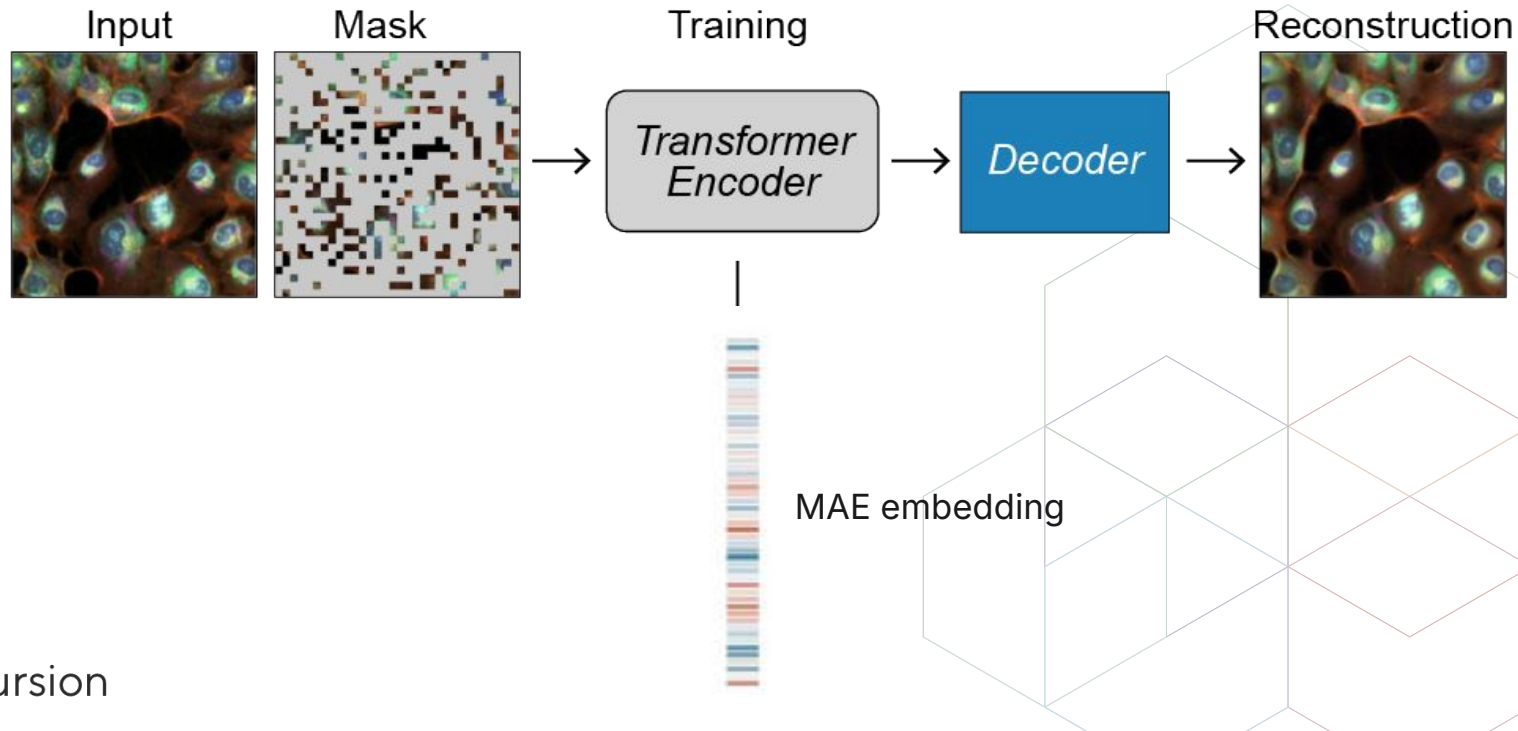
Kian Kenyon-Dean, Zitong Jerry Wang, John Urbanik, Konstantin Donhauser, Jason Hartford, Saber Saberian, Nil Sahin, Ihab Bendi, Safiye Celik, Marta Fay, Juan Sebastián Rodríguez Vera, Imran S Haque, Oren Kraus

Kian Kenyon-Dean @ NeurIPS FM4Sci Workshop 2024
Staff Machine Learning Engineer, Recursion

Context & related work



Masked Autoencoding (MAE): *effective self-supervised learning when scaled*



Phenom-1: Large Vision Transformer (ViT-L)

- ViT-Large/8, 330 million parameters (1024+1 tokens per sample)
- MAE trained on RPI-93M for ~40 epochs

Masked Autoencoders for Microscopy are Scalable Learners of Cellular Biology

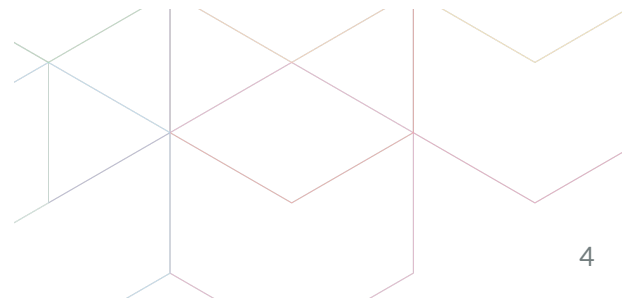
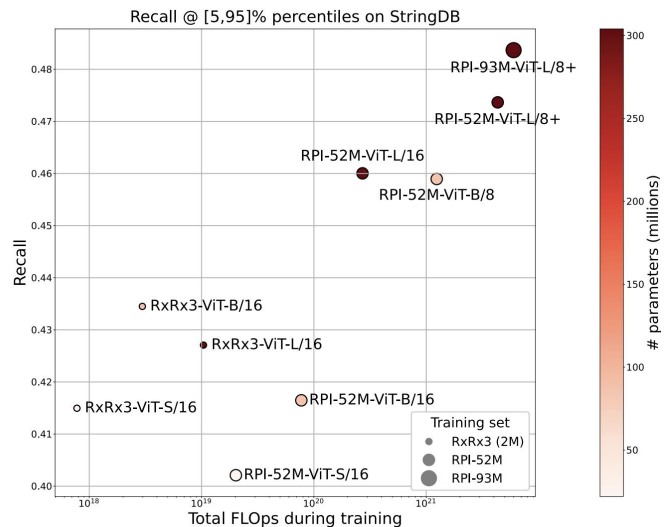
Oren Kraus¹ Kian Kenyon-Dean¹ Saber Saberian¹ Maryam Fallah¹ Peter McLean¹
Jess Leung¹ Vasudev Sharma¹ Ayla Khan¹ Jia Balakrishnan¹ Safiye Celik¹
Dominique Beaini² Maciej Sypetkowski² Chi Vicky Cheng¹ Kristen Morse¹
Maureen Makes¹ Ben Mabey¹ Berton Earnshaw^{1,2}
¹Recursion ²Valence Labs

CVPR
JUNE 17-21, 2024



Masked Autoencoders are Scalable Learners of Cellular Morphology

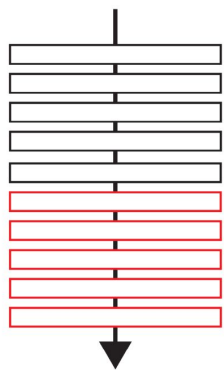
Oren Kraus* Kian Kenyon-Dean* Saber Saberian Maryam Fallah Peter McLean
Jess Leung Vasudev Sharma Ayla Khan Jia Balakrishnan Safiye Celik
Maciej Sypetkowski Chi Vicky Cheng Kristen Morse Maureen Makes
Ben Mabey Berton Earnshaw



Our contributions in this work

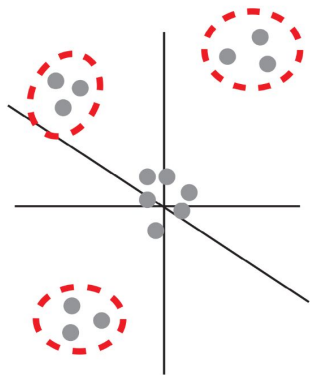
General foundation model training approach

Model scaling

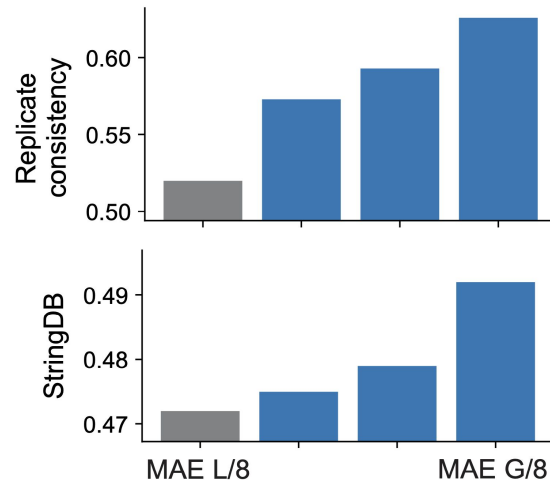
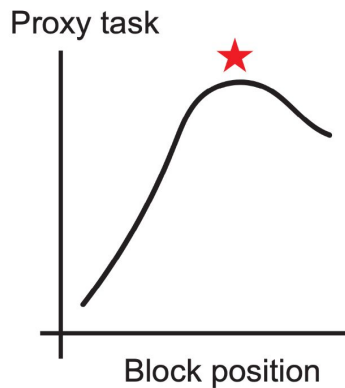


307M → 1.9B

Curated data



Block search



Strategies

Block search
Curated data
Model scaling



	Phenom-1 (MAE-L/8)	Phenom-2 (MAE-G/8)
Training Regime	Self-supervised MAE	Self-supervised MAE
Model architecture	Vision Transformer (ViT-L/8) 24 blocks	Vision Transformer (ViT-G/8) 48 blocks
Model parameters	330 million	1.9 BILLION (+6x)
Dataset base	93M unique well images	Specially curated Phenoprints-16M (-6x)
Dataset sampling	3.5 billion crops == 3.6 trillion tokens	8 billion crops (+2.2x) == 8.2 trillion tokens
Training time	20,000 A100 hours	43,000 H100 hours (+4x) (i.e., >5 GPU-years!)

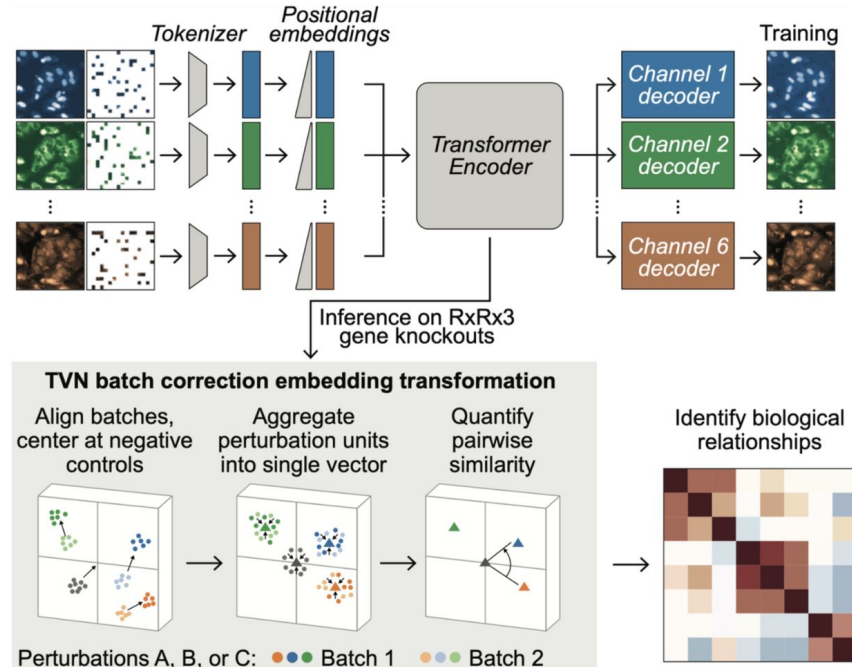
Channel-Agnostic MAE baseline

- CA-MAE-S/16 (1536+1 tokens per sample)
 - Trained on RxRx3 for 100 epochs
- **OpenPhenom** - new open-source publicly available model on HuggingFace 🙌
 - CA-MAE-S/16 but also trained on RxRx3 + public JUMP-CP data.



Evaluation: how “good” (in terms of batch effect correction / *replicate consistency* and *relationship prediction*) are your model’s **genetic representations** derived from images?

- **Inference regime:** 36 center crops ($256 \times 256 \times 6$) per well ($2048 \times 2048 \times 6$ pixels) → **80 million RxRx3 images fed forward through a single trained model** to have a comprehensive **whole-genome** evaluation \$\$



Whole-Genome analysis is expensive and labour-intensive, so let's **linearly probe** the model on a small dataset to find *the best layer*.

New task: Anax

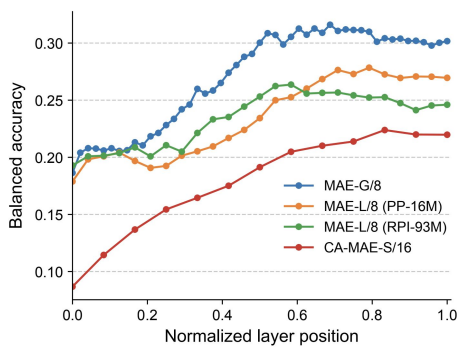
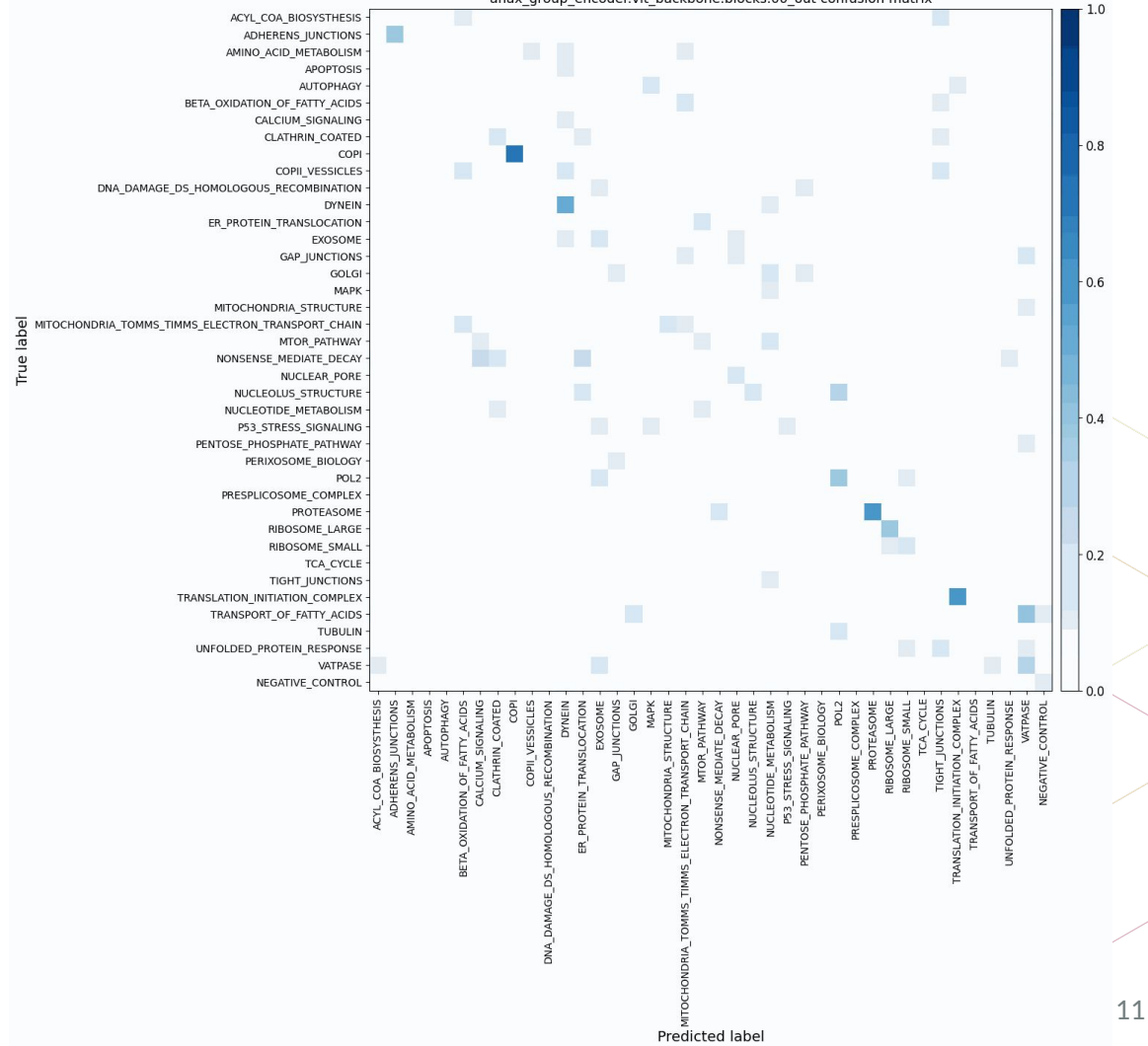
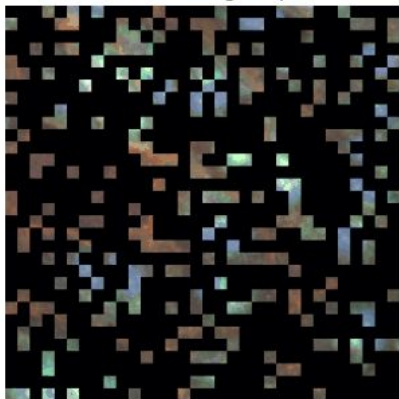


Table 4: Anax groups and their associated genes. This table presents a comprehensive list of gene groups and their corresponding genes.

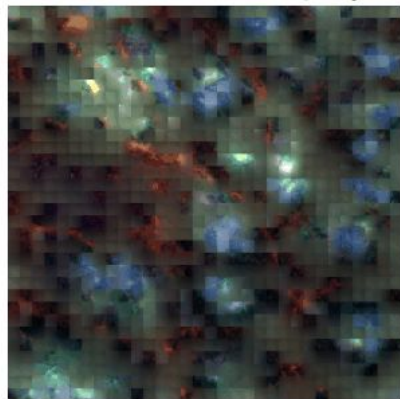
Anax Group	Genes
Acyl Coa Biosynthesis	ELOVL2, ELOVL5, ELOVL6, HACD1, HACD2, HSD17B12, SCD, SCD5, TCCR
Adherens Junctions	ACTB, ACTG1, AFN2, CDH1, CTNNA1, CTNNA3, CTNND1, NECTIN1, NECTIN3, NECTIN4
Amino Acid Metabolism	ALDH1A1, ARG2, CRB1, CSMD2, CPB1, DMO, GTC, PKCK1, PKCK3, SMT
Apoptosis	CFLAR, DFFB, CASP8, CASP9, FASLG, BCL2, DFFA, XIAP, TNFSF10, ACT3
Autophagy	ATG12, ATG3, ATG8B, ATGAC, ATG7, GABARAP, PRKC3, PRK3B, PRKAA1, ULK1
Beta Oxidation Of Fatty Acids	ACAA2, ACADL, ACADM, ACADS, ACADVL, BEH1, ECH1, HADH, HADHA, HADHB
Calcium Signaling	ADCY1, ADCY2, ADCY3, CALM1, CAMEB1, CAMEB2, PDB1B, PDB1C, PRKAC1, PRKAC2
Clathrin Coated Vesicles	AP2A1, AP2A2, AP2B1, AP2M1, AP2S1
COP1	ARCN1, COB1, COPB1, COPB2, COPE, COPG1, COPG2
COP1 Vesicles	SEC13, SEC23A, SEC8B, SEC8D, SEC14
DNA Damage Repair	BLM, BRCA2, EME1, NEMO, POLR2, RADD1B, RADD1C, RADD1D, RPA1, XRCC2
Dystin	DYNC1H1, DYNC1H2, DYNC1L1, DYNC1L2, DYNLL1
ER Protein Translocation	SPC3, SEC13A1, SRP14, SRP72, SRP51, SRP6, SEC11A, SRP98, SRP98, SRP94
Endosome	EEA1, EEA2, EEA3, EEA4, EEA5, EEA6, EEA7, EEA8, EEA9, EEA10, EEA11, EEA12, EEA13, EEA14, EEA15, EEA16, EEA17, EEA18, EEA19, EEA20, EEA21, EEA22, EEA23, EEA24, EEA25, EEA26, EEA27, EEA28, EEA29, EEA30, EEA31, EEA32, EEA33, EEA34, EEA35, EEA36, EEA37, EEA38, EEA39, EEA40, EEA41, EEA42, EEA43, EEA44, EEA45, EEA46, EEA47, EEA48, EEA49, EEA50, EEA51, EEA52, EEA53, EEA54, EEA55, EEA56, EEA57, EEA58, EEA59, EEA60, EEA61, EEA62, EEA63, EEA64, EEA65, EEA66, EEA67, EEA68, EEA69, EEA70, EEA71, EEA72, EEA73, EEA74, EEA75, EEA76, EEA77, EEA78, EEA79, EEA80, EEA81, EEA82, EEA83, EEA84, EEA85, EEA86, EEA87, EEA88, EEA89, EEA90, EEA91, EEA92, EEA93, EEA94, EEA95, EEA96, EEA97, EEA98, EEA99, EEA100
Gap Junctions	ADCY8, DRD2, HTR2C, HTR2D, HTR2E, HTR2F, HTR2G, HTR2H, HTR2I, HTR2J, HTR2K, HTR2L, HTR2M, HTR2N, HTR2O, HTR2P, HTR2Q, HTR2R, HTR2S, HTR2T, HTR2U, HTR2V, HTR2W, HTR2X, HTR2Y, HTR2Z
Golgi	ACTR10, ACTR11A, CAZPA3, COG8, CTSC, PPP9C, RAB11B, SEC23C, SEC23G, TME8D9
MAPK	ERK5, ERK6, FERF1, FERF2, HSP1, MAP2K2, MAP2K3, RAC1, RAC1A, RASGEF3
Mitochondria Structure	APOD, APOD1, TMEH1, CHCHD6, AITPM6, MICOS1, AITP51C, DNAI3, DMAC1L, AITP51P
Mitochondrial Transport	ATP5FA, COXA, COX4, COX17, HSPA9, PPTM1, PMP1C, PMP1B, SLC25A4
mTOR Pathway	CAB39, CAB39L, EIF4B1, MESTR, PRKAA2, RPS6KB1, RPTOR, STK11, STRADA, TSC1
Nonsense Mediated Decay	CAC1, EIF4A3, MAGOR, MAGORB, RIBK1
Nuclear Pore	NUP17, NUP133, NUP133, NUP188, NUP210, NUP97, NUP85, NUP93
Nucleolin Structure	FBL, NAT10, NOL1, NOP58, UTP9
Nucleolin Metabolism	ADSL, ADSL1, ADSL2, ATRC, CSM5, IMPDH1, IMPDH2, PAK3, PAK5, PPA1
PS3 Stress Signaling	ATM, ATR, CENPL, CDK1, CHEK1, CHEK2, MDM2, MDM1, TP53, TP53
Pentose Phosphate Pathway	GRPD, TALDO1, DERA, RPE, PGMD, RBK5, PGM, PGL5, RPEL1, PRPS2
Perovastome Biology	ACOT9, AQP5, BAAT, HMGCL, HSD17B4, MEVD3, PAOC, PAOC2, PAOC3, PEK3, PIP5K
Proteasome Complex	AAV1B1, AOR, CENPL1, EXO1, FANSD1, LAMA1, RPL1, PRK4, SANS1, SENS4
Proteasome	PSMA1, PSMA4, PSMB1, PSMB2, PSMB7, PSMA6, PSMA3, PSMB4, PSMA5, PSMB3
Ribosome Large	RPL13A, RPL11, RPL10, RPL23A, RPL30, RPL7A, RPLP2, RPL28, RPL5, RPL27A
Ribosome Small	RPS28, RPS18, RPS19, RPS18, RPS11, RPS3A, RPS19, RPS15, RPS45, RPS9
RNA Polymerase II	POLR2A, POLR2B, POLR2C, POLR2D, POLR2E, POLR2F, POLR2G
TCA Cycle	AC02, DLST, FH, IDH1, IDH3B, MDH1, OGDH, SDHB, SLC142, SLC162
Tight Junctions	CLDN14, CLDN17, CLDN18, CLDN19, CLDN4, CLDN5, CLDN9, MIPPS, PAR6B, PRCK1
Translation Initiation Complex	EIF3E, EIF3A, EIF3F, EIF3L, EIF3K, EIF3P1, EIF3B, EIF3H, EIF3I, EIF3J, EIF3L
Transport Of Fatty Acids	APOD, LCN12, LCN15, LCN9, SLC27A1, SLC27A4, SLC27A6
Tubulin	TUBA3C, TBCD, TUBA8, TUBA8, TUBA3, TUBA1A, TUBB4B, ARL2, TUBA1B
Unfolded Protein Response	CXK1, DNAJB11, ERF25, KHSBP, MIFP1, SIRT1, TAC1D2, TLR1, TSPY12, YP1A
V-ATPase	ATP9A, ATP9B, ATP9V12, ATP9V11, ATP9V10, ATP9V11



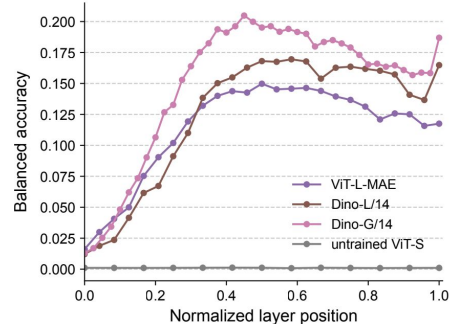
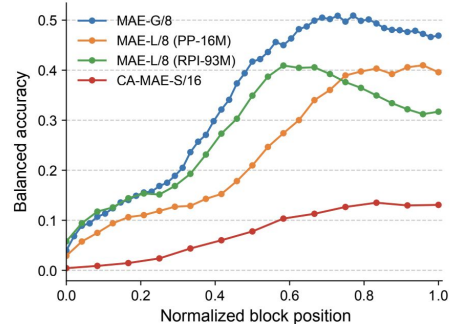
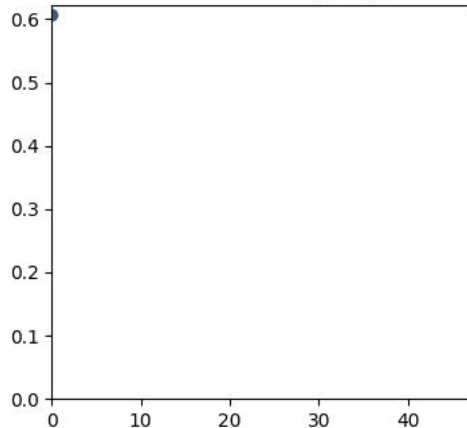
Masked Image input



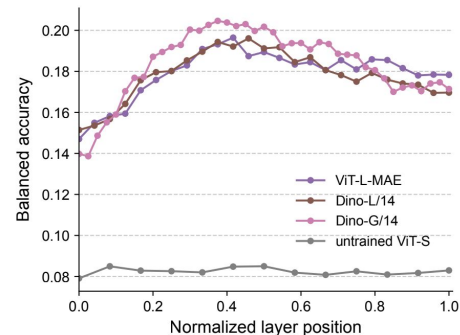
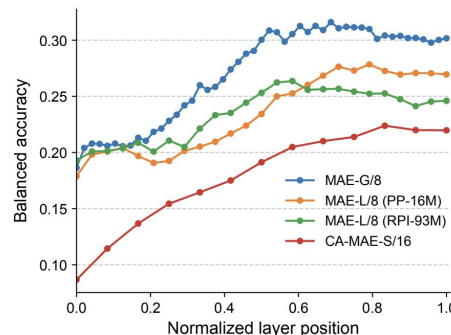
Phenom-2 Reconstruction @ Layer 0



Difference with original image



(a) RxRx1 siRNA knockdown classification.

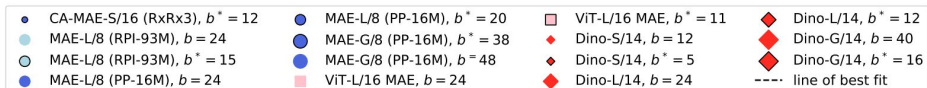
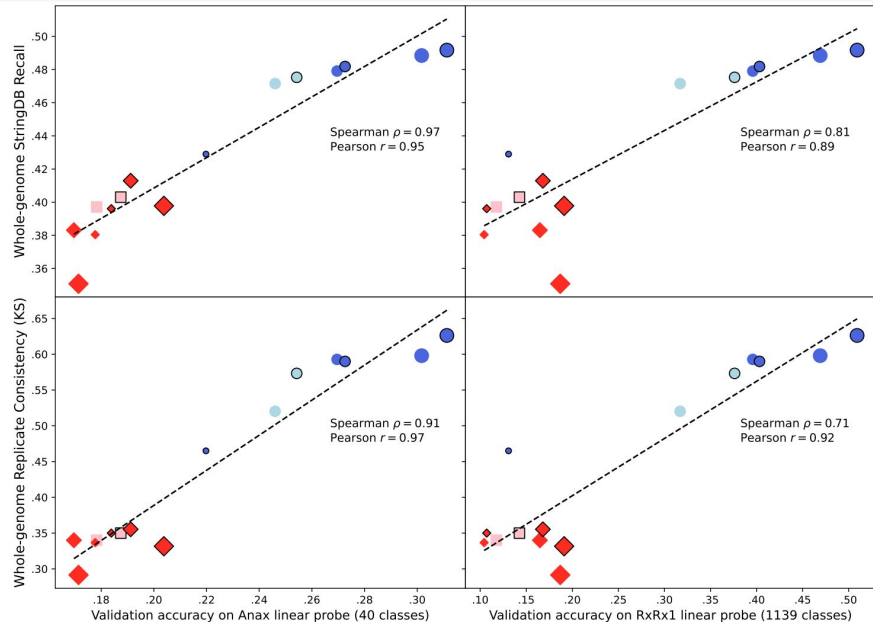


(b) Anax functional gene group classification.

Figure 3: Block-wise validation set **linear probe results** comparing ViT models pretrained on cell microscopy images (left) versus natural images (right). (a) 1139-class RxRx1 SiRNA knockdown classification (Sypetkowski et al., 2023); (b) 40-class Anax functional gene group classification on HUVEC cell images from RxRx3 CRISPR knockouts (Fay et al., 2023).

Final takeaways

- **Curated microscopy data** = awesome
- Nearly all **SSL ViTs** we evaluate (MAEs and Dino-v2 imagenet baselines) are **better at intermediate layers**
- **Linear separability on small datasets strongly correlates to performance at the whole-genome scale and transfer to new datasets for these SSL models**
- Scaled **MAE to 1.9 billion parameters** is **SOTA** across variety of **newly evaluated benchmarks**



Model backbone	b	Pretraining data	CORUM	hu.MAP	Reactome	StringDB
CellProfiler	-	N/A	.219	.184	.131	.191
CA-MAE-S/16	12	RxRx3	.233	.199	.154	.214
MAE-L/8	24	RPI-93M	.248	.208	.160	.226
MAE-G/8	38	Phenoprints-16M	.264	.215	.165	.235

Table 2: Biological relationship recall benchmarks at 0.05-0.95 cosine threshold on public JUMP-CP image data (Chandrasekaran et al., 2023) generated by completely different labs and assay protocols compared to the data used for pretraining. Each result has a standard deviation $\leq \pm 0.0023$, and spans nearly 8,000 gene-knockouts and are computed after applying PCA with center-scaling for embedding post-processing alignment.

Questions?

kian.kd@recursion.com

oren.kraus@recursion.com

info@rxrx.ai

