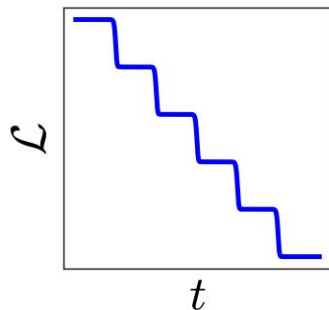


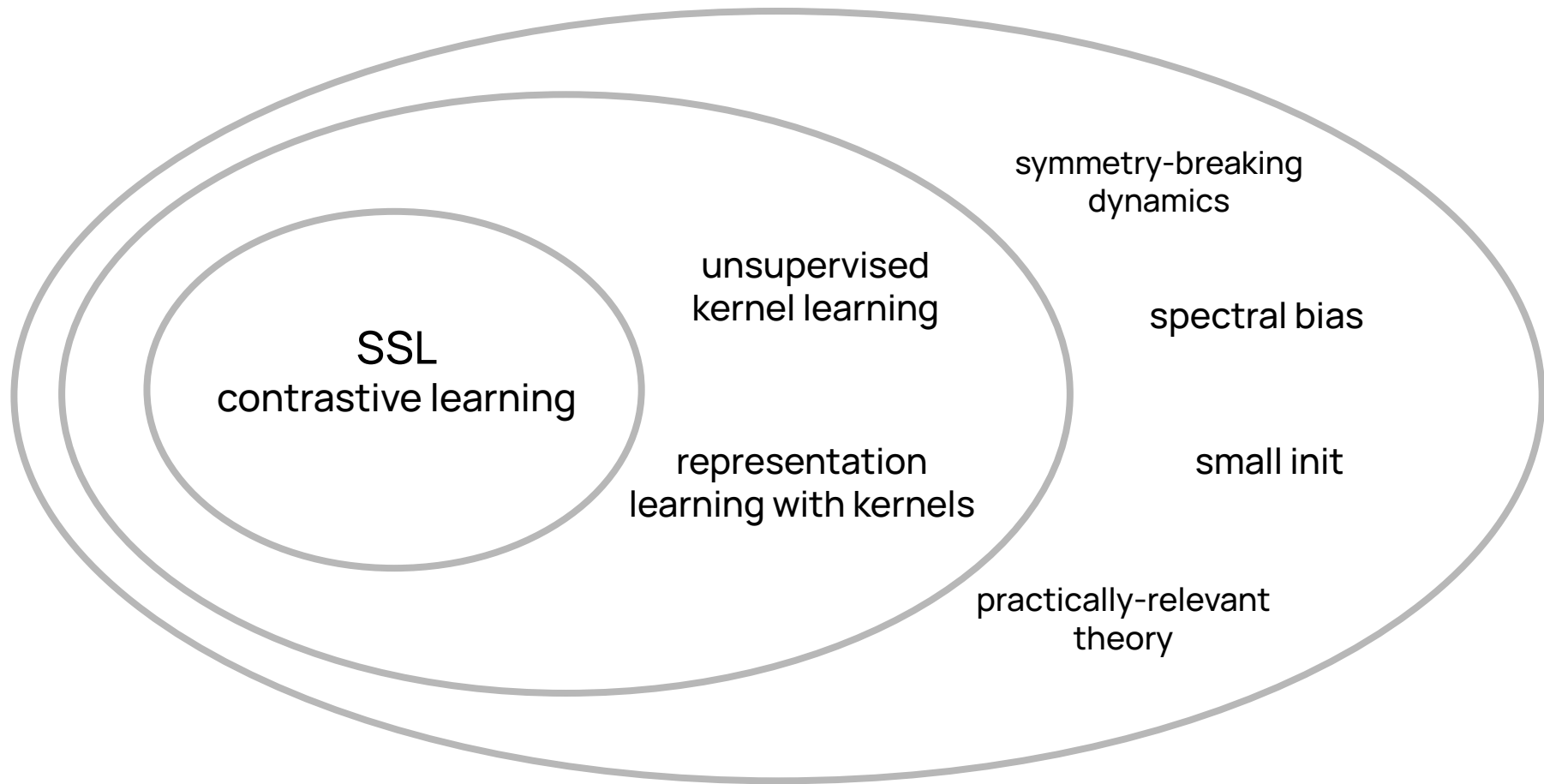
On the stepwise nature of SSL



Jamie Simon

with Maksis Knutins, Liu Ziyin, Daniel Geisz, Abe
Fetterman, and Josh Albrecht

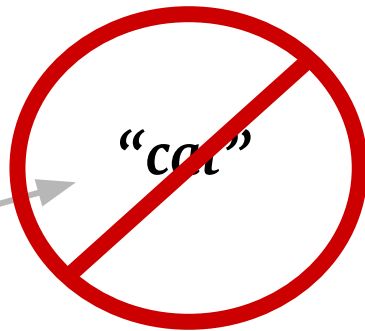
paper: tinyurl.com/stepwise-ssl-paper



As theorists, we ought to look beyond supervised learning.

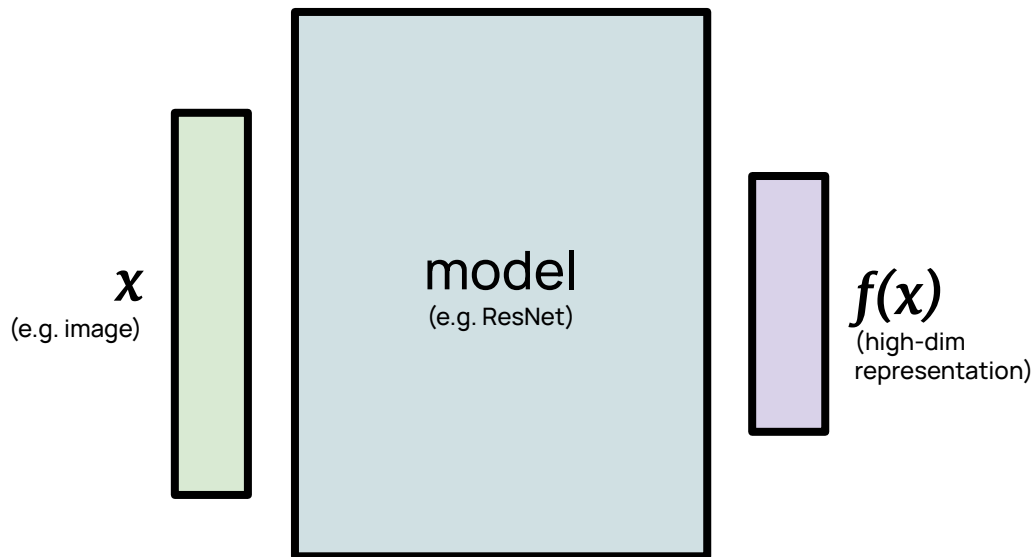


We should think about learning *representations*, not *target functions*.



(high-dim
representation)

Joint embedding SSL learns representations with a simple prescription.



let x, x' be two related samples.

goal of training: **make $f(x)$ close to $f(x')$**
(but not degenerate)

Surprisingly, this works very well...

SimCLR, Barlow Twins, VICReg, BYOL, DINO, CLIP

A Simple Framework for Contrastive Learning of Visual Representations

Ting Chen¹ Simon Kornblith¹ Mohammad Norouzi¹ Geoffrey Hinton¹



SSL is state-of-the-art in image domains.

...but here's what we understand about it:

[this slide intentionally left blank]

What is the nature of the learning process of SSL?

What representations are ultimately learned?

Idea: answer these Qs in the linear case and check agreement with deep nets.

Plan of attack

1. Find and solve an exactly solvable case of SSL with a **linear model**.
2. Kernelize it.
3. Check agreement with realistic SSL.

Linear model of **Barlow Twins** (Zbontar et al. 2021)

Let the dataset contain n “positive pairs” $(\mathbf{x}_i, \mathbf{x}'_i)_{i=1}^n$.

Let the trained function be linear as $\mathbf{f}(\mathbf{x}) = \mathbf{W}\mathbf{x}$.

Correlation matrix: $\mathbf{C} \equiv \frac{1}{2n} \sum_i (\mathbf{f}(\mathbf{x}_i)\mathbf{f}(\mathbf{x}'_i)^\top + \mathbf{f}(\mathbf{x}'_i)\mathbf{f}(\mathbf{x}_i)^\top)$

Barlow Twins loss: $\mathcal{L} = \|\mathbf{C} - \mathbf{I}_d\|_F^2$

Train with GD: $\frac{d\mathbf{W}}{dt} = -\nabla_{\mathbf{W}}\mathcal{L}$.

Key insight: diagonalize w.r.t. the “feature cross-covariance.”

$$\mathbf{\Gamma} \equiv \frac{1}{2n} \sum_i (\mathbf{x}_i \mathbf{x}_i'^\top + \mathbf{x}_i' \mathbf{x}_i^\top)$$

$$\Rightarrow \mathcal{L} = \|\mathbf{W}\mathbf{\Gamma}\mathbf{W}^\top - \mathbf{I}_d\|_F^2$$

$$\Rightarrow \frac{d\mathbf{W}}{dt} = -4 (\mathbf{W}\mathbf{\Gamma}\mathbf{W}^\top - \mathbf{I}_d) \mathbf{W}\mathbf{\Gamma}$$

^This is *exactly solvable* when $\mathbf{W}_0 = \mathbf{W}(t=0)$ is either *aligned* or *small*.

Aligned init: $\mathbf{W}_0 = \mathbf{U} \mathbf{S}_0 \mathbf{\Gamma}^{(\leq d)}$

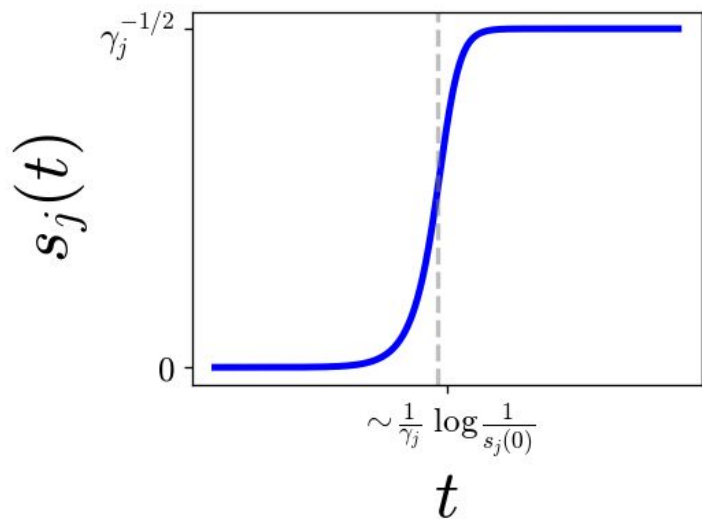
right singular vectors are
the top d eigenvectors of $\mathbf{\Gamma}$

From aligned init, each singular value evolves independently.

Proposition 4.1 (Trajectory of $\mathbf{W}(t)$ from aligned initialization). *If $\mathbf{W}(0) = \mathbf{W}_0$ as given by Equations 7 and 8, then*

$$\mathbf{W}(t) = \mathbf{U}\mathbf{S}(t)\mathbf{\Gamma}^{(\leq d)} \quad (9)$$

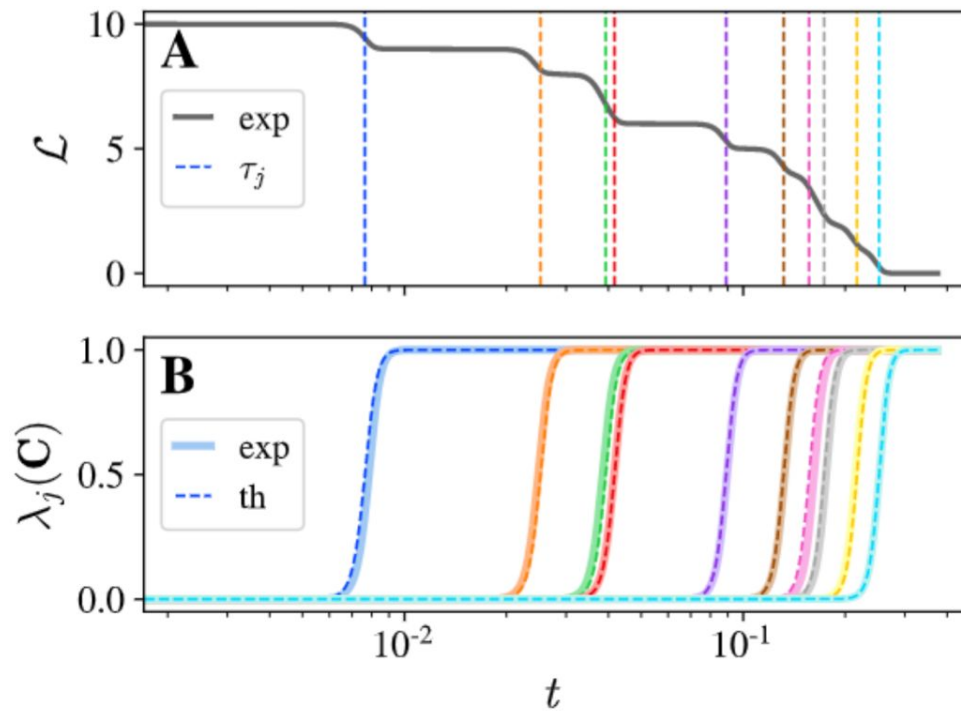
with $\mathbf{S}(t) = \text{diag}(s_1(t), \dots, s_d(t))$ and



$$s_j(t) = \frac{e^{4\gamma_j t}}{\sqrt{s_j^{-2}(0) + (e^{8\gamma_j t} - 1)\gamma_j}}. \quad (10)$$

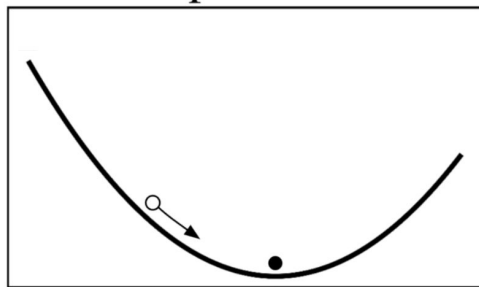
\uparrow
 j -th eigenvalue of $\mathbf{\Gamma}$

From *generic small init*, these dynamics follow the aligned trajectory!



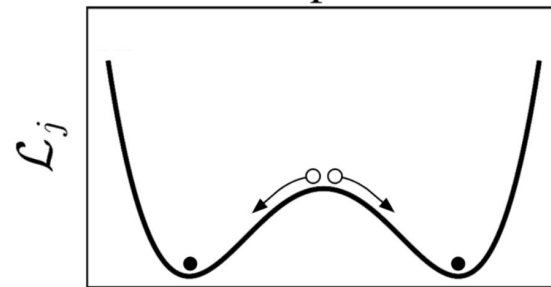
Representation learning is a process of *symmetry breaking*.

Supervised

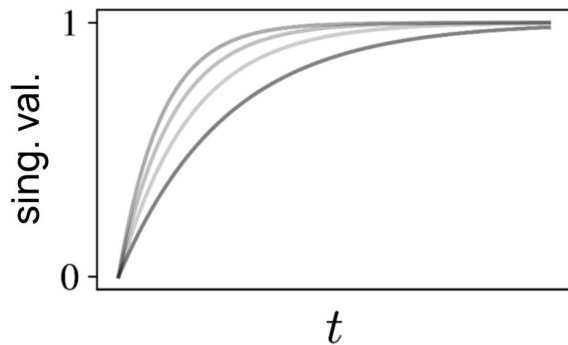


0 1
sing. val.

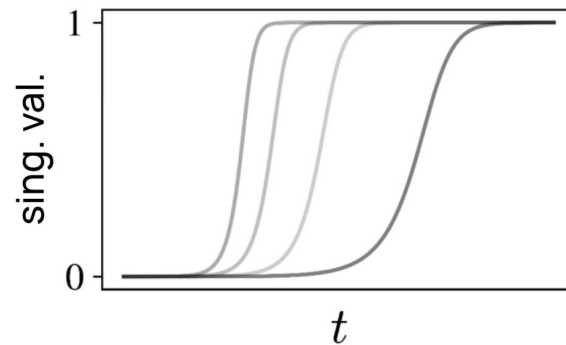
Self-supervised



-1 0 1
sing. val.



$$\dot{s}_j = \gamma_j (s_j^* - s_j)$$



$$\dot{s}_j = \gamma_j s_j (1 - s_j^2)$$

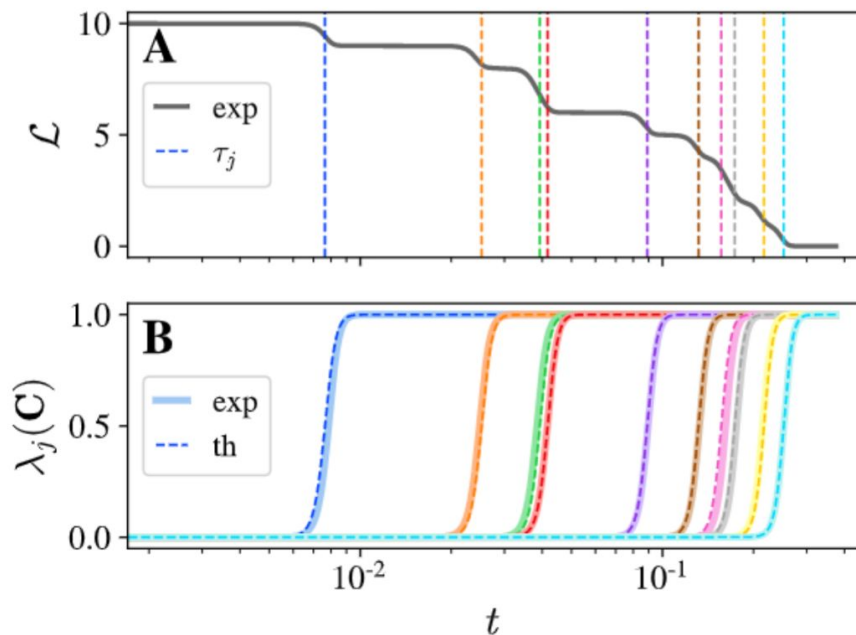
It's kernel time

To kernelize, rephrase everything in terms of inner products $\mathbf{x}^T \mathbf{x}'$, then sub in $K(\mathbf{x}, \mathbf{x}') \equiv \mathbf{x}^T \mathbf{x}'$.

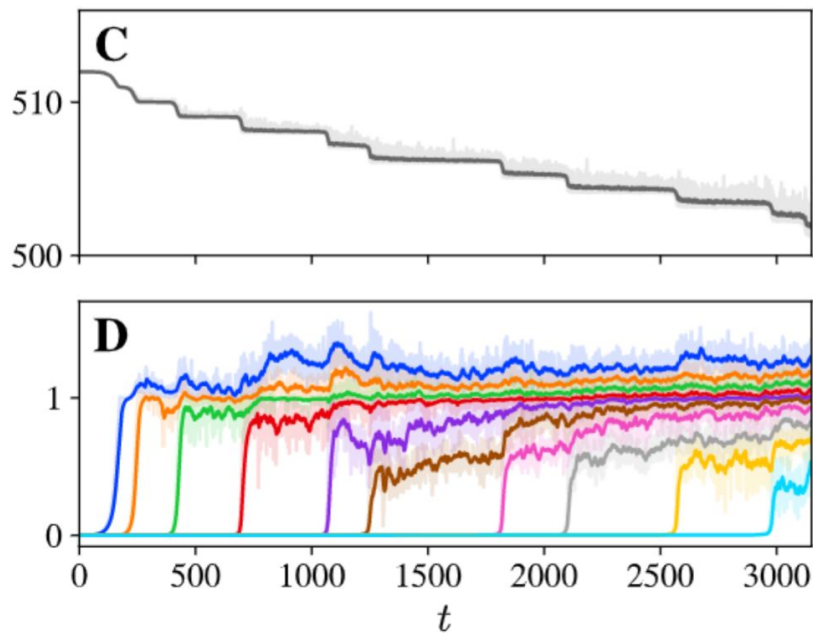
End of theory. Let's check against ResNets.

ResNets with small init show stepwise learning on STL-10.

Linear toy model

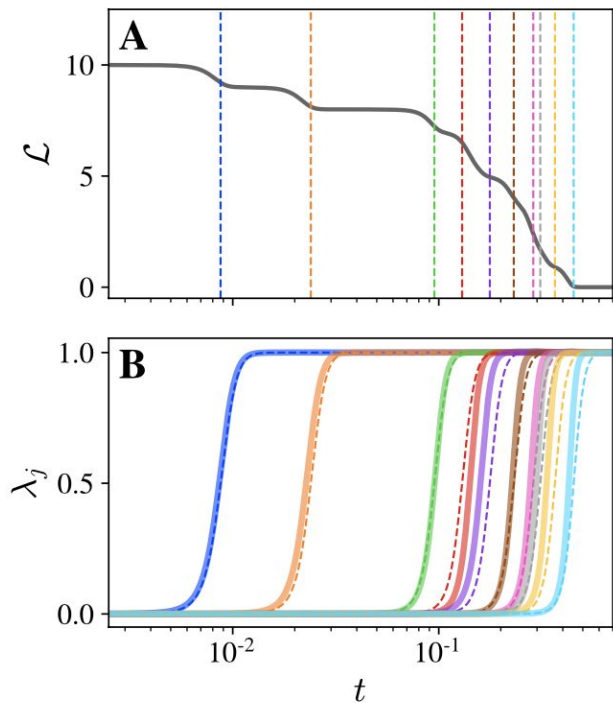


Barlow Twins (ResNet-50)

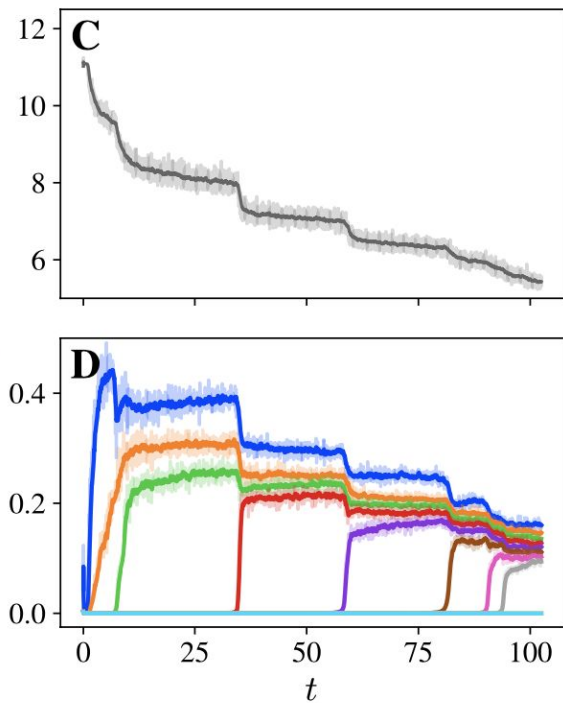


ResNets with small init show stepwise learning on STL-10 (cont.).

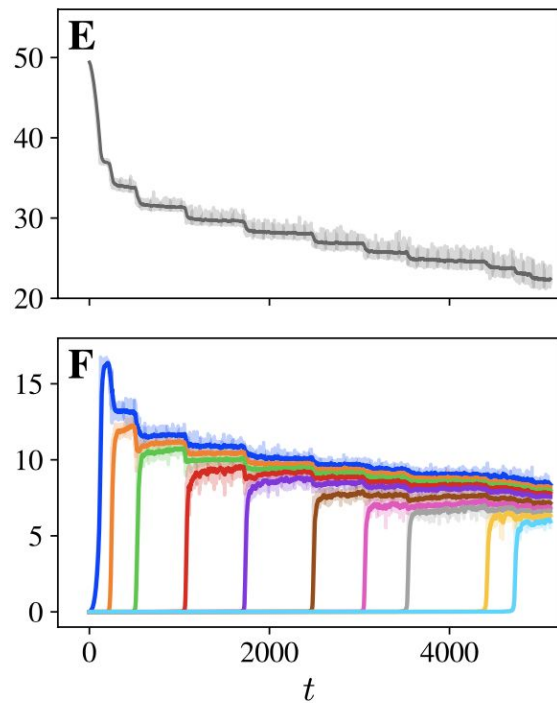
Shallow MLP



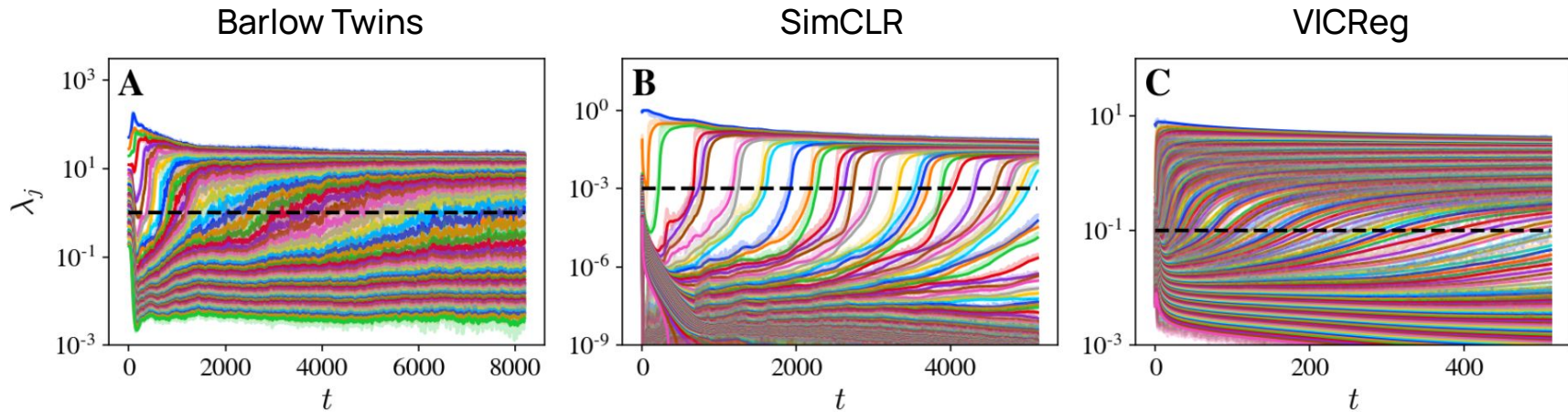
SimCLR (ResNet-50)



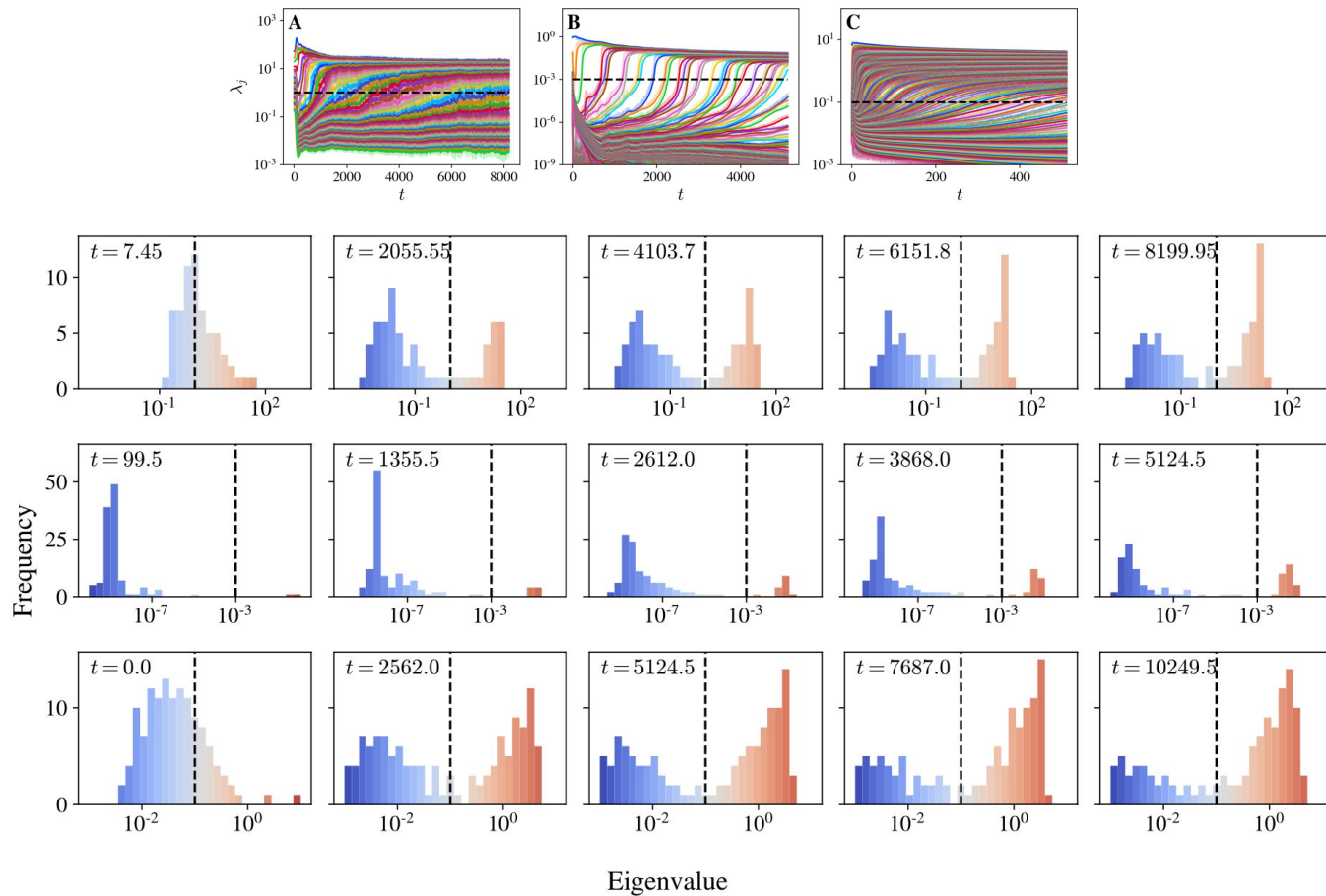
VICReg (ResNet-50)



...and you can still see the stepwise behavior with regular init!



...and you can still see the stepwise behavior with regular init!



Future directions

Training SSL is famously slow. Can we speed it up?

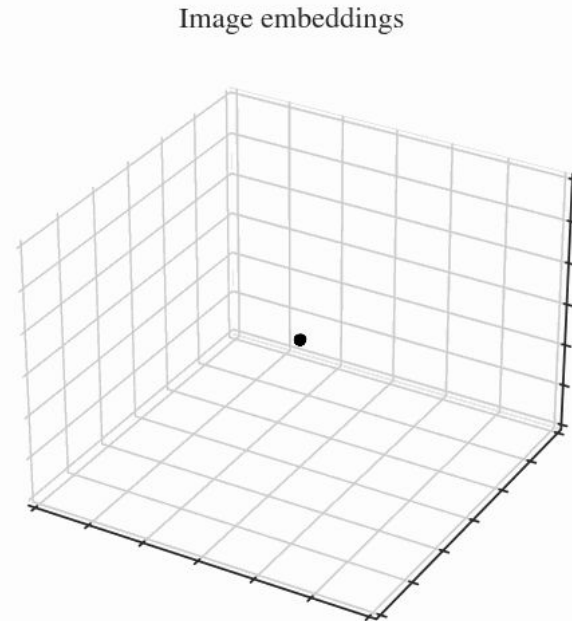
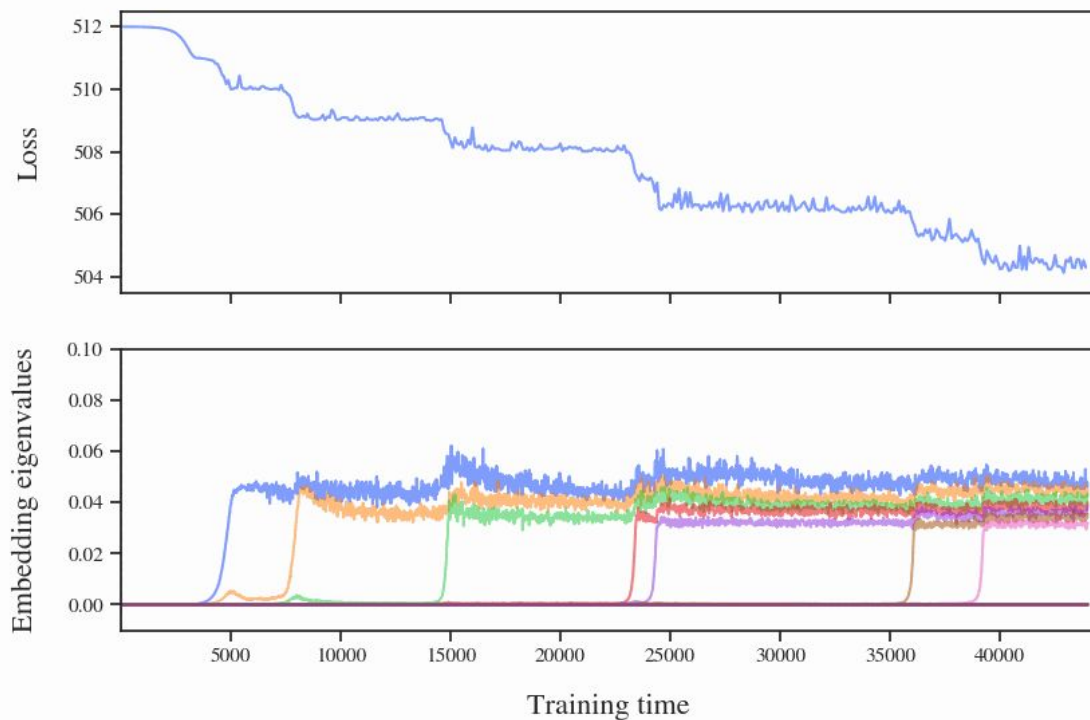
Interpretability of top eigendirections?

Stepwise learning in other settings? LLMs?

Broader lessons

- 1) Let's study (self/un)-supervised learning
- 2) It's worth collaborating with experimentalists

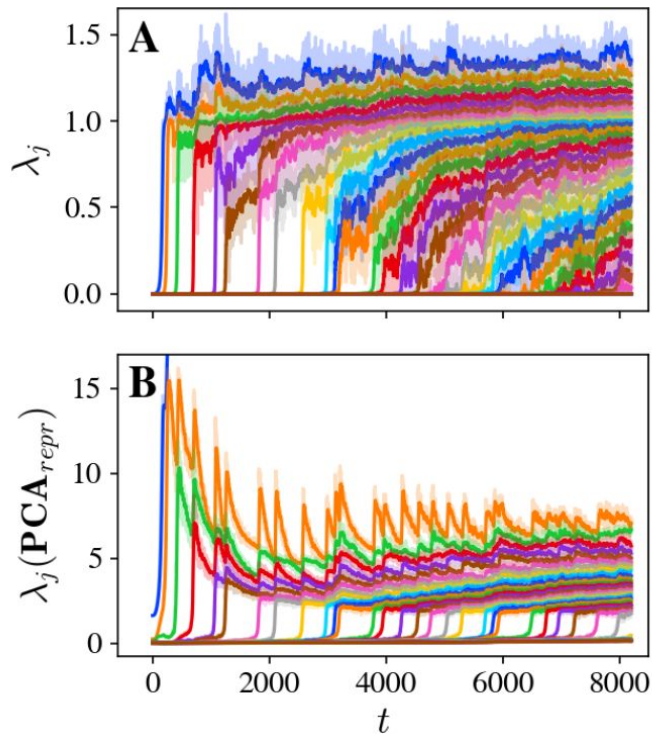
Summary: SOTA self-supervised systems show striking steps ('specially starting small')



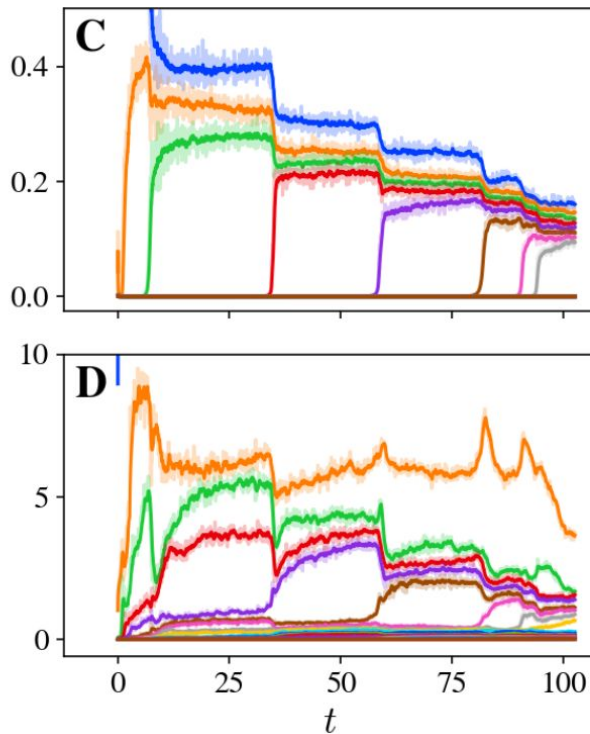
(backups)

...and you can see it in *hidden representations* too!

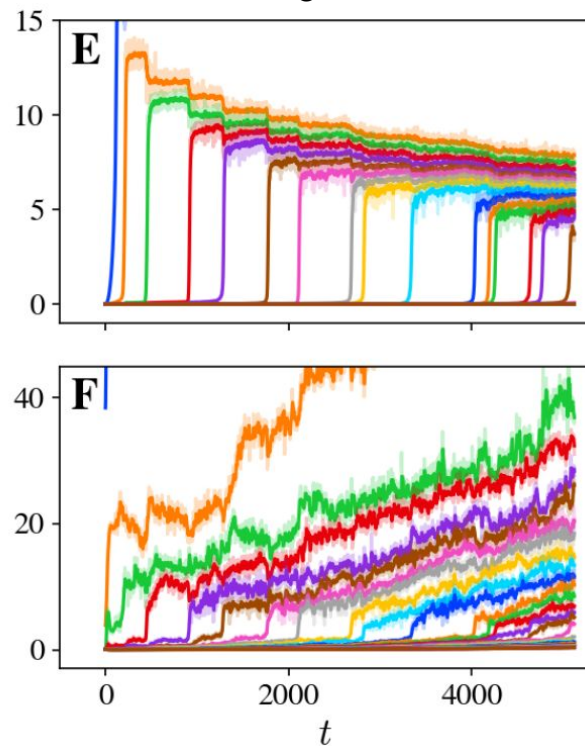
Barlow Twins



SimCLR



VICReg



We can mostly predict embeddings from the final NTK.

Procedure:

- Train a ResNet on full STL-10 dataset from small init.
- Measure eNTK after training on 1k positive pairs.
- Compute predicted embeddings from the eNTK.
- Compute inner product of true + predicted embedding subspaces.

	BT	SimCLR	VICReg	(random subspaces)
$a(\mathbf{P}_{\text{exp}}, \mathbf{P}_{\text{th}})$	0.615	0.517	0.592	0.025

Not perfect, but pretty good agreement!