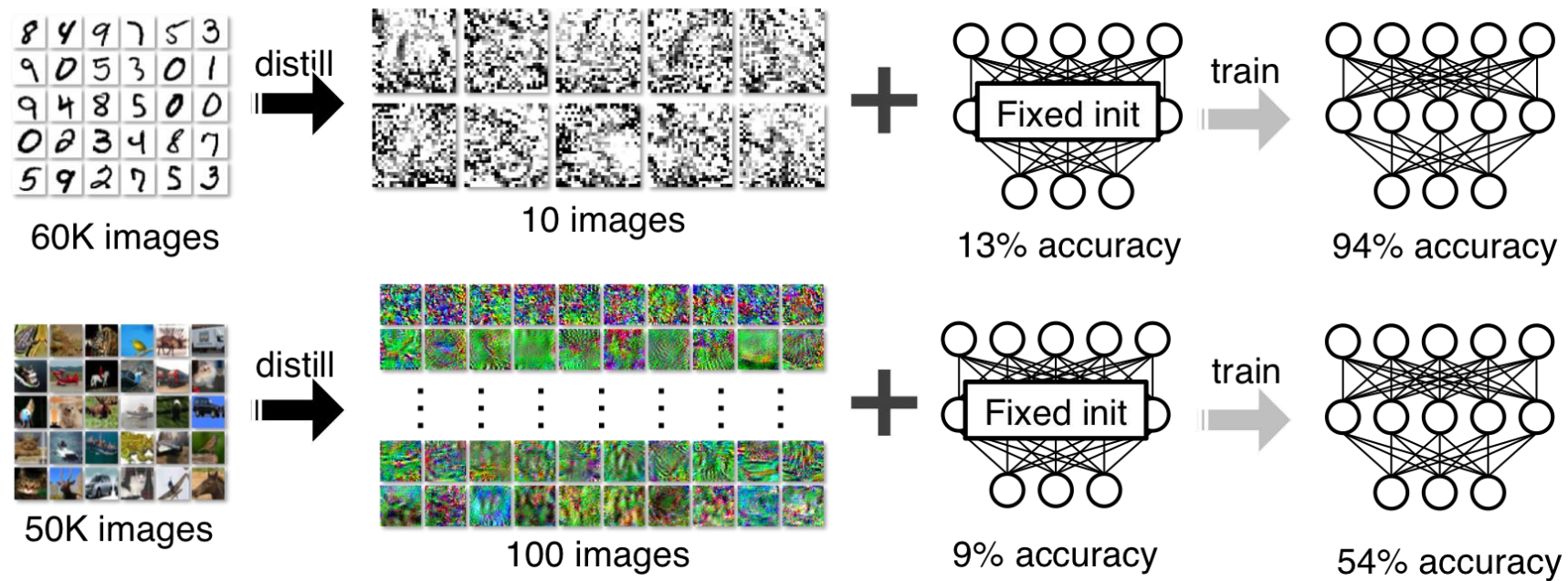


# Dataset Distillation with Infinitely Wide Convolutional Networks

Mikhail Kuznetsov, 4th year, CMC MSU

# What is data distillation?



Dataset distillation on MNIST and CIFAR10

# Applications

# Distillation method.Recap

Ridge Regression

$$\min_w \|Xw - y\|_2^2 + \lambda \|w\|_2^2$$

$$w^* = (X^\top X + \lambda I)^{-1} X^\top y$$

$$y_{\text{new}} = X_{\text{new}} w^*$$

Kernel Ridge Regression

$$\min_w \|K_{XX} w - y\|_2^2 + \lambda w^\top K_{XX} w$$

$$w^* = (K_{XX} + \lambda I)^{-1} y$$

$$y_{\text{new}} = K_{X_{\text{new}}X} w^*$$

# Distillation method.KIP

$$L(X_s, y_s) = \frac{1}{2} \left\| y_t - K_{X_t X_s} (K_{X_s X_s} + \lambda I)^{-1} y_s \right\|_2^2$$

---

**Algorithm 1:** Kernel Inducing Point (KIP )

---

**Require:** A target labeled dataset  $(X_t, y_t)$  along with a kernel or family of kernels.

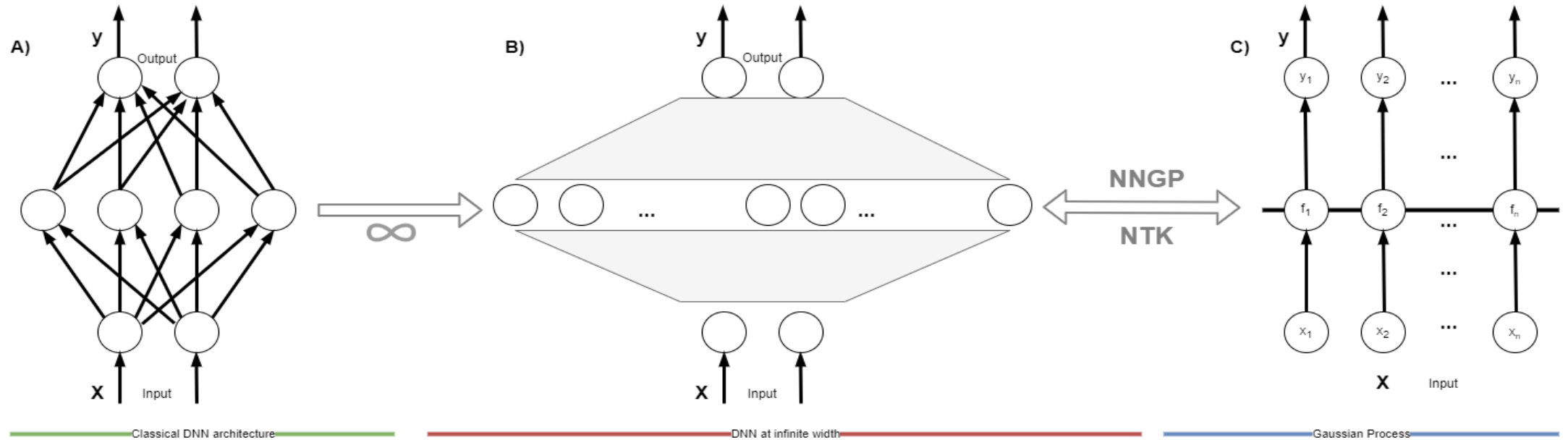
- 1: Initialize a labeled support set  $(X_s, y_s)$ .
  - 2: **while** not converged **do**
  - 3:   Sample a random kernel. Sample a random batch  $(\bar{X}_s, \bar{y}_s)$  from the support set. Sample a random batch  $(\bar{X}_t, \bar{y}_t)$  from the target dataset.
  - 4:   Compute the kernel ridge-regression loss given by (7) using the sampled kernel and the sampled support and target data.
  - 5:   Backpropagate through  $\bar{X}_s$  (and optionally  $\bar{y}_s$  and any hyper-parameters of the kernel) and update the support set  $(X_s, y_s)$  by updating the subset  $(\bar{X}_s, \bar{y}_s)$ .
  - 6: **end while**
  - 7: **return** Learned support set  $(X_s, y_s)$
-

# Distillation method.Algorithm

- 1) Forward pass: compute kernels by batch partition
- 2) Backward pass:

$$\frac{\partial L}{\partial X_s} = \frac{\partial L}{\partial (K(X_s, X_s))} \frac{\partial K(X_s, X_s)}{\partial X_s} + \frac{\partial L}{\partial (K(X_t, X_s))} \frac{\partial K(X_t, X_s)}{\partial X_s}$$

# Infinitely Wide Convolutional Networks



Neural Tangents [[lib](#)] [[paper](#)]

# Preprocessing.ZCA-regularized

- 1) flatten the features for each train image and then standardize each feature across the train dataset.
- 2) feature-feature covariance matrix  $C = U\Sigma U^T$ ,
- 3) Let  $\bar{W}_\lambda = U\phi_\lambda(\Sigma)U^T$  where  $\phi_\lambda: \mu \text{ to } (\mu + \lambda\overline{\text{tr}}C)^{-1/2}$ ,  $\overline{\text{tr}}(C) = \text{tr}(C)/\text{len}(C)$ .
- 4) New features: standardize + @  $\bar{W}_\lambda$

(if  $\lambda = 0 \rightarrow$  standard ZCA = I cov matrix)



# Experiments.KIP vs ALL

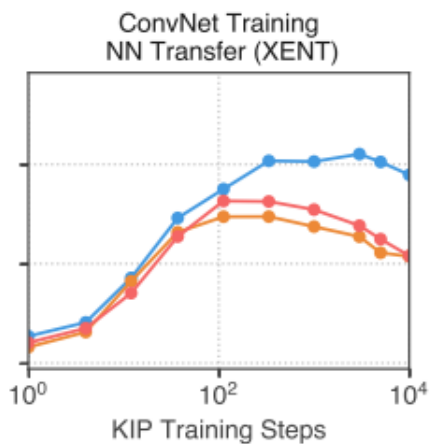
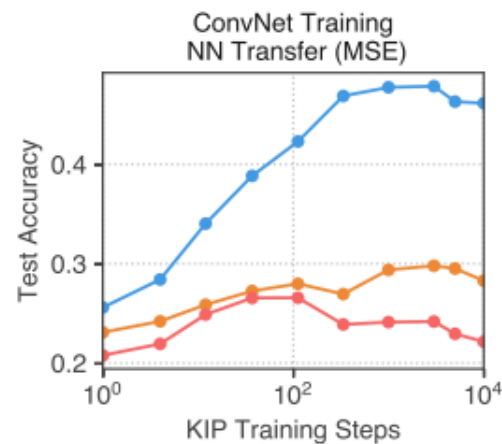
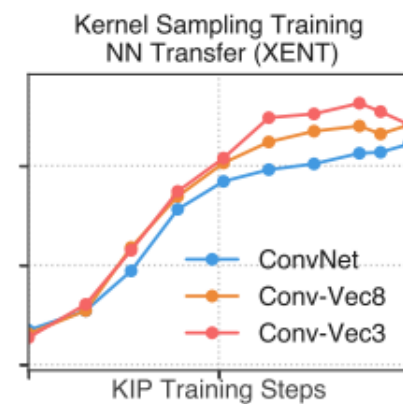
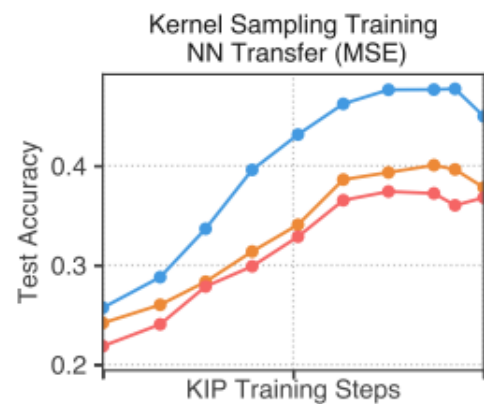
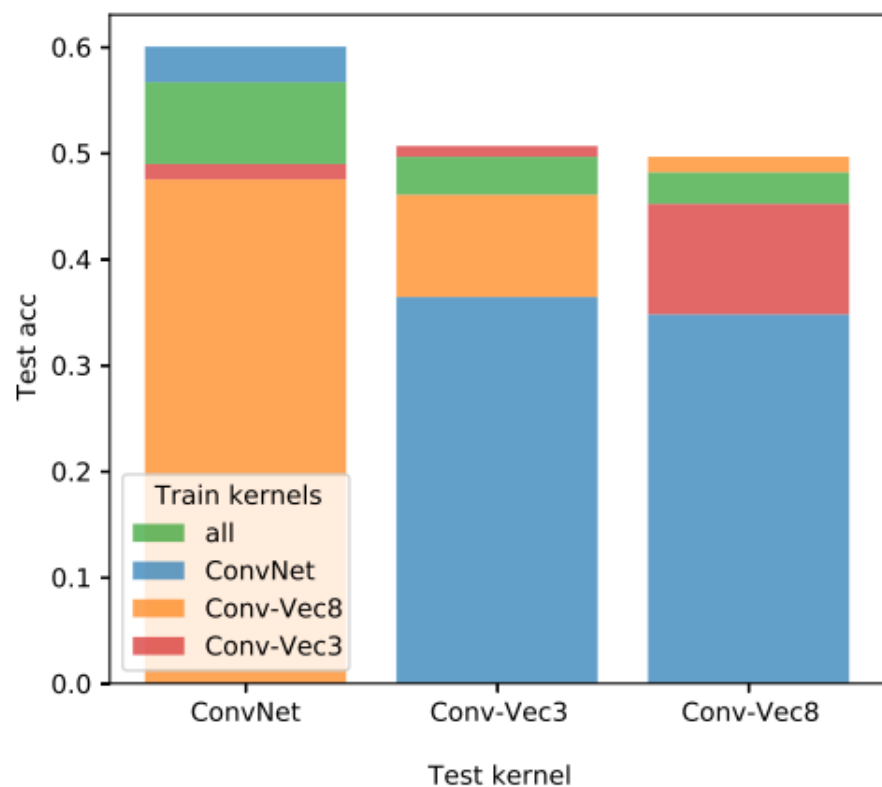
Table 1: **Comparison with other methods.** The left group consists of neural network based methods. The right group consists of kernel ridge-regression. All settings for KIP involve the use of label-learning. Grayscale datasets use standard channel-wise preprocessing while RGB datasets use regularized ZCA preprocessing.

	Imgs/ Class	DC <sup>1</sup>	DSA <sup>1</sup>	KIP FC <sup>1</sup> aug	LS ConvNet <sup>2,3</sup>	KIP ConvNet <sup>2</sup>	
						no aug	aug
MNIST	1	91.7±0.5	88.7±0.6	85.5±0.1	73.4	<b>97.3±0.1</b>	<b>96.5±0.1</b>
	10	97.4±0.2	97.8±0.1	97.2±0.2	96.4	<b>99.1±0.1</b>	<b>99.1±0.1</b>
	50	98.8±0.1	99.2±0.1	98.4±0.1	98.3	<b>99.4±0.1</b>	<b>99.5±0.1</b>
Fashion-MNIST	1	70.5±0.6	70.6±0.6	-	65.3	<b>82.9±0.2</b>	76.7±0.2
	10	82.3±0.4	84.6±0.3	-	80.8	<b>91.0±0.1</b>	88.8±0.1
	50	83.6±0.4	88.7±0.2	-	86.9	<b>92.4±0.1</b>	91.0±0.1
SVHN	1	31.2±1.4	27.5±1.4	-	23.9	62.4±0.2	<b>64.3±0.4</b>
	10	76.1±0.6	79.2±0.5	-	52.8	79.3±0.1	<b>81.1±0.5</b>
	50	82.3±0.3	<b>84.4±0.4</b>	-	76.8	82.0±0.1	<b>84.3±0.1</b>
CIFAR-10	1	28.3±0.5	28.8±0.7	40.5±0.4	26.1	<b>64.7±0.2</b>	63.4±0.1
	10	44.9±0.5	52.1±0.5	53.1±0.5	53.6	<b>75.6±0.2</b>	<b>75.5±0.1</b>
	50	53.9±0.5	60.6±0.5	58.6±0.4	65.9	78.2±0.2	<b>80.6±0.1</b>
CIFAR-100	1	12.8±0.3	13.9±0.3	-	23.8	<b>34.9±0.1</b>	33.3±0.3
	10	25.2±0.3	32.3±0.3	-	39.2	47.9±0.2	<b>49.5±0.3</b>

<sup>1</sup> DSA from [1], DC from [2], KIP FC from [3], LS ConvNet from [4], KIP ConvNet from [5].

+ ZCA, + label-learning, +- aug

# Experiments.Kernels Choice

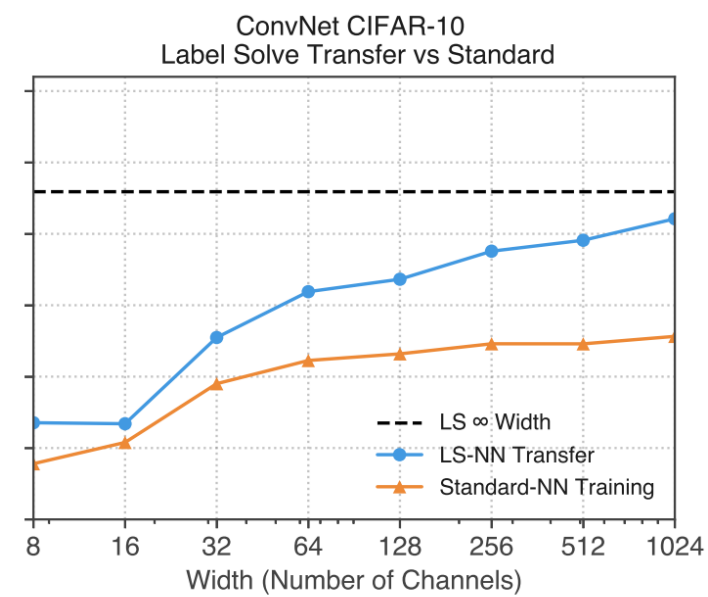
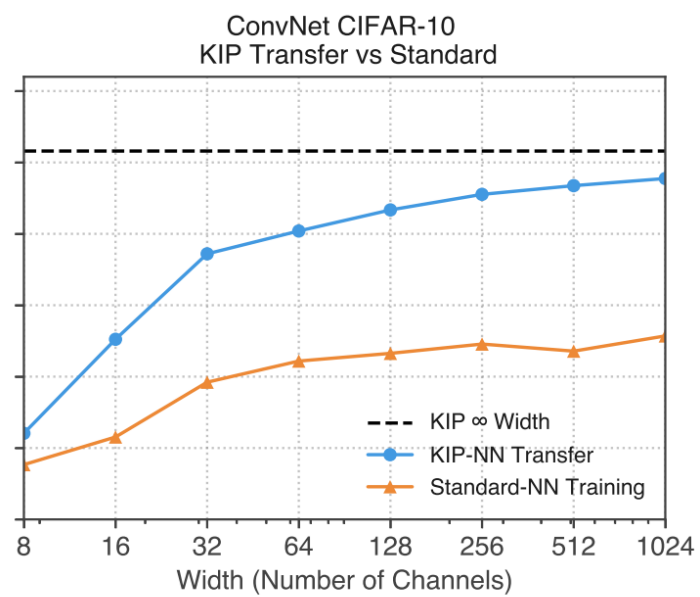
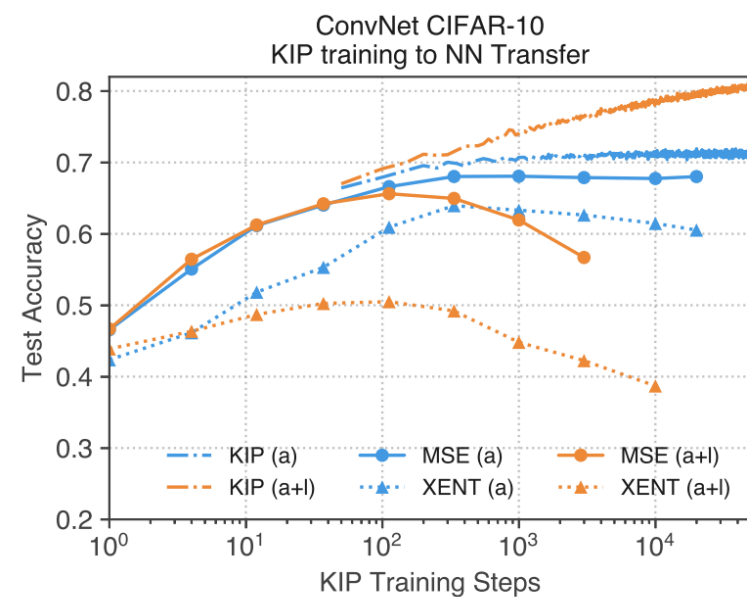


# Experiments.Transfer

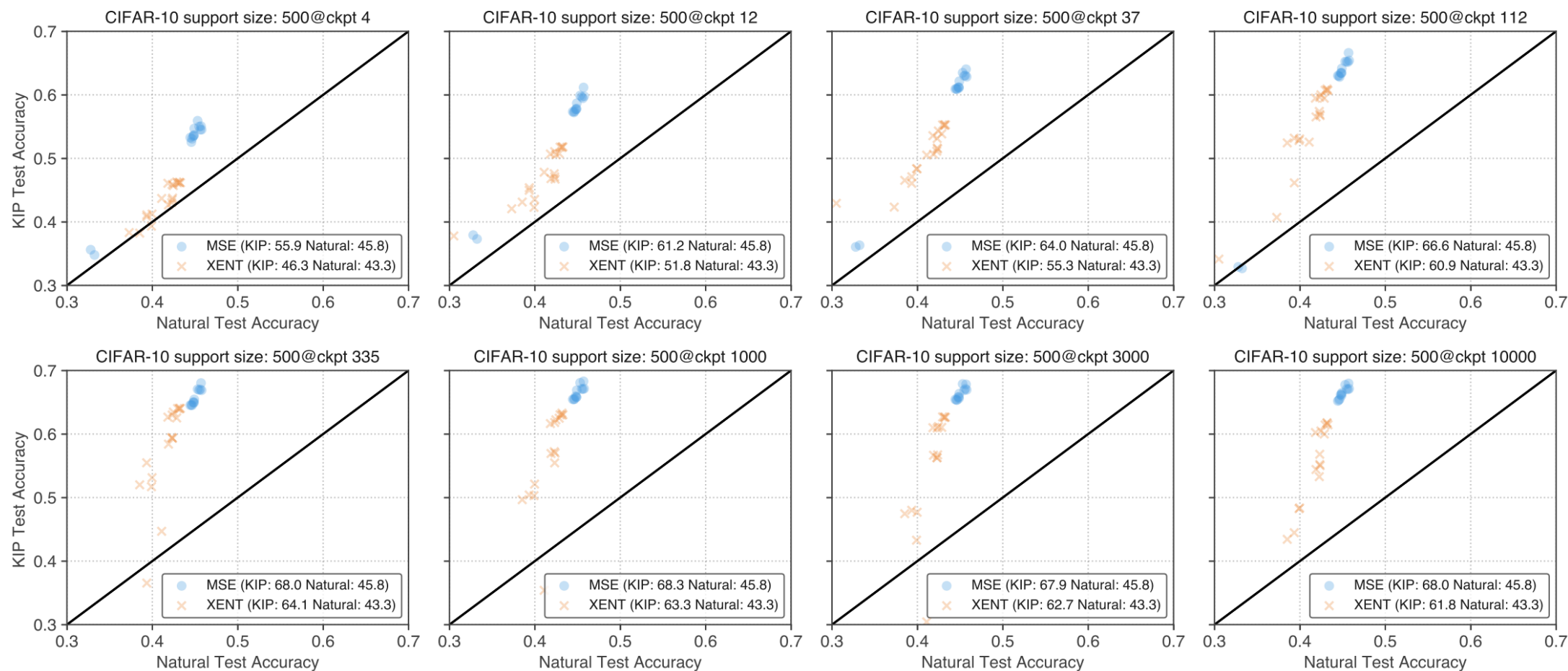
	Imgs/Class	DC/DSA	KIP to NN	Perf. change	LS to NN	Perf. change
MNIST	1	<b>91.7±0.5</b>	90.1±0.1	-5.5	71.0±0.2	-2.4
	10	<b>97.8±0.1</b>	97.5±0.0	-1.1	95.2±0.1	-1.2
	50	<b>99.2±0.1</b>	98.3±0.1	-0.8	97.9±0.0	-0.4
Fashion-MNIST	1	70.6±0.6	<b>73.5±0.5*</b>	-9.8	61.2±0.1	-4.1
	10	84.6±0.3	<b>86.8±0.1</b>	-1.3	79.7±0.1	-1.2
	50	<b>88.7±0.2</b>	88.0±0.1*	-4.5	85.0±0.1	-1.8
SVHN	1	31.2±1.4	<b>57.3±0.1*</b>	-8.3	23.8±0.2	-0.2
	10	<b>79.2±0.5</b>	75.0±0.1	-1.6	53.2±0.3	0.4
	50	<b>84.4±0.4</b>	80.5±0.1	-1.0	76.5±0.3	-0.4
CIFAR-10	1	28.8±0.7	<b>49.9±0.2</b>	-9.2	24.7±0.1	-1.4
	10	52.1±0.5	<b>62.7±0.3</b>	-4.6	49.3±0.1	-4.3
	50	60.6±0.5	<b>68.6±0.2</b>	-4.5	62.0±0.2	-3.9
CIFAR-100	1	13.9±0.3	<b>15.7±0.2*</b>	-18.1	11.8±0.2	-12.0
	10	<b>32.3±0.3</b>	28.3±0.1	-17.4	25.0±0.1	-14.2

ZCA preprocessing, +- aug, \* - best with trained labels

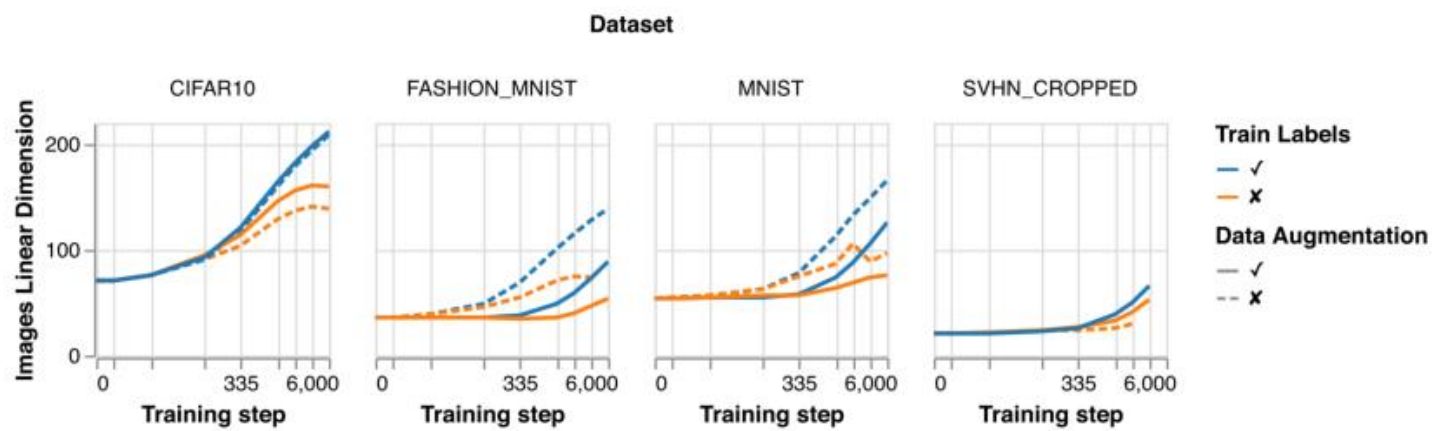
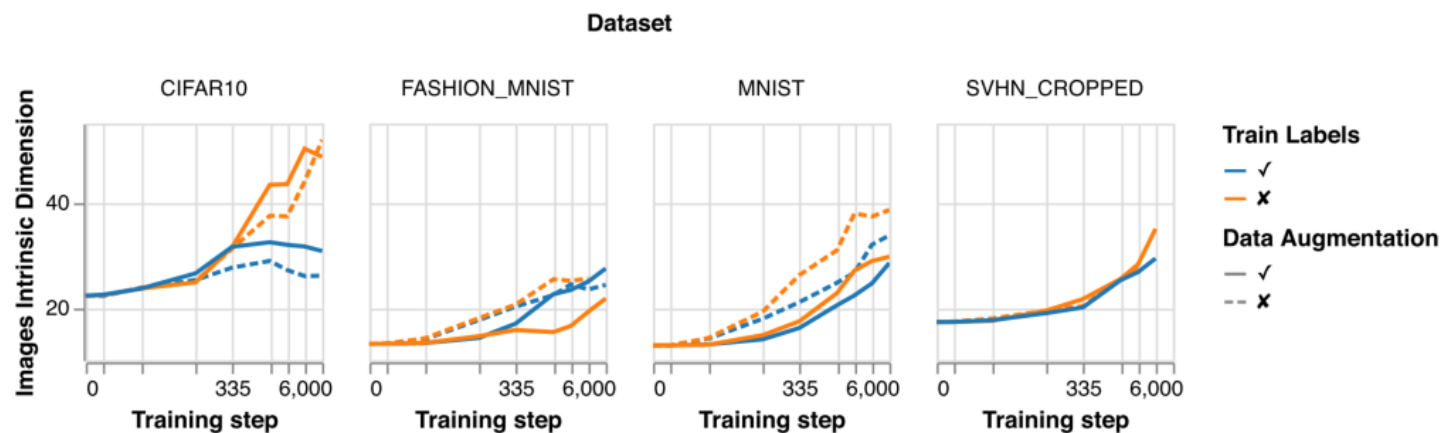
# Experiments.Transfer



# Experiments.Hyperparameters



# Experiments.Datasets Dims



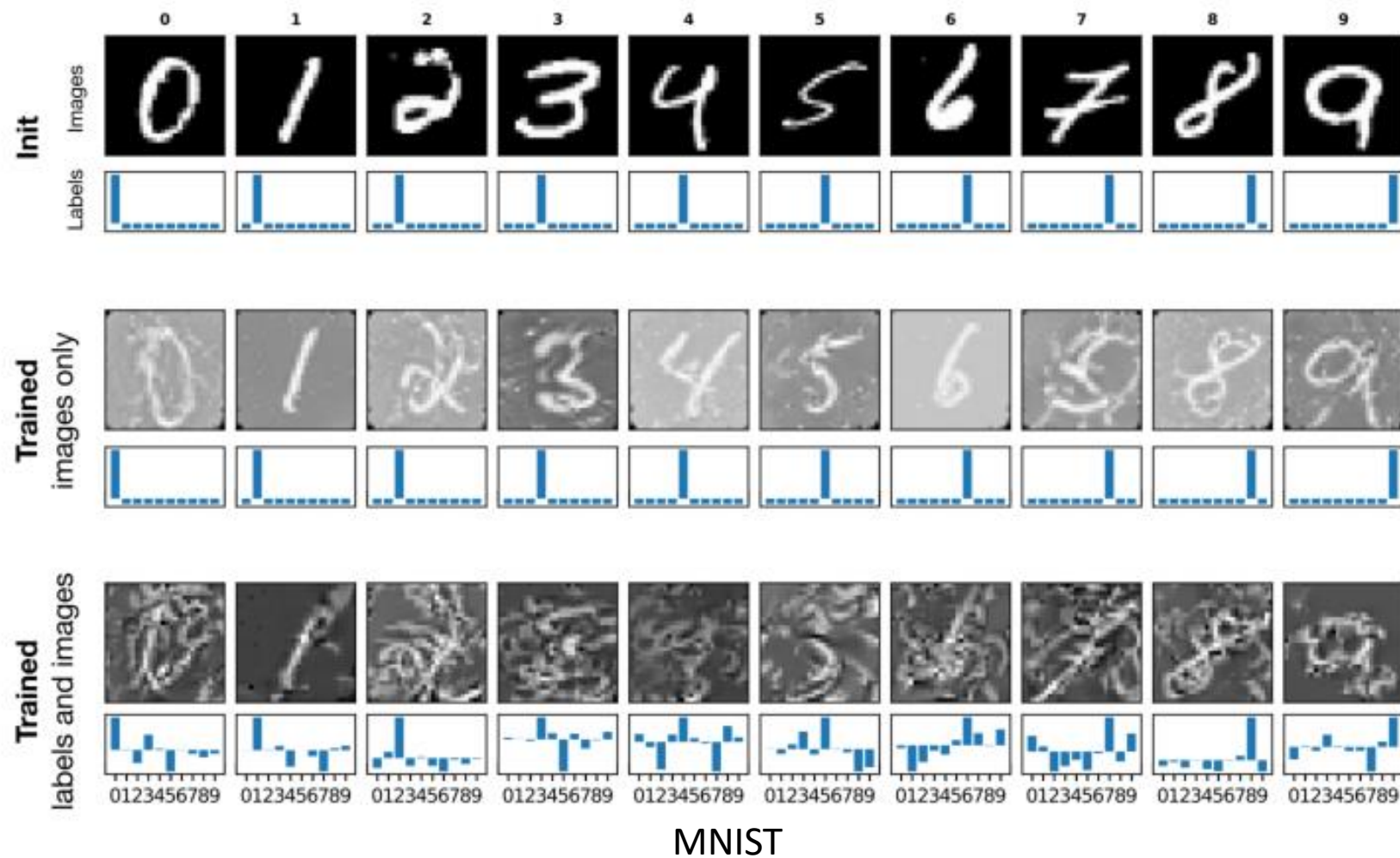
# Experiments.Learned Images



CIFAR-100



# Experiments.Learned Images





# Papers

- [2] Wang, T., Zhu, J., Torralba, A. and Efros, A. Dataset Distillation, 2018.