# Learning Curves for Analysis of Deep Networks

Derek Hoiem[1]

Tanmay Gupta[2]
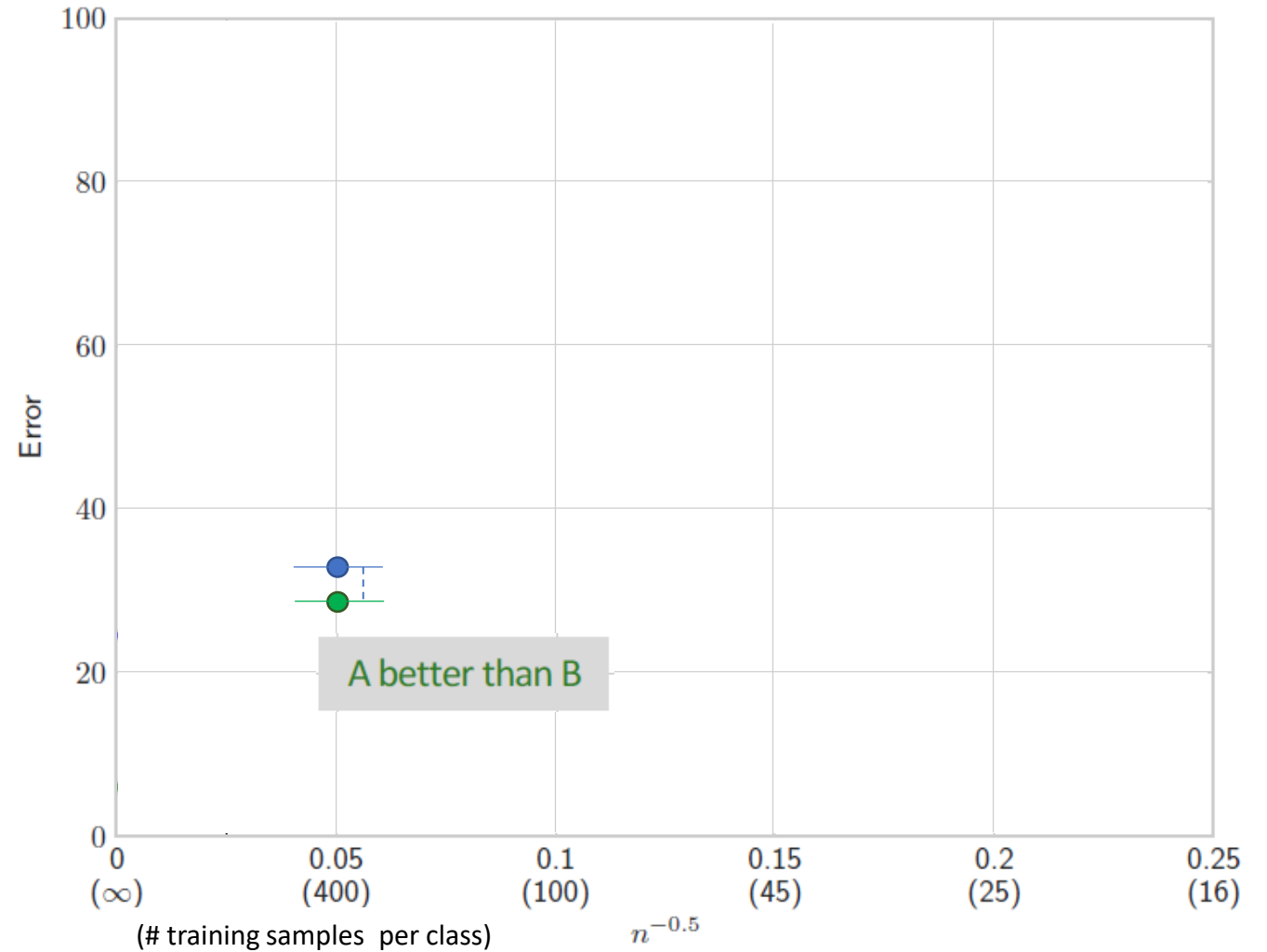
Zhizhong Li[1]

Michal M. Shlapentokh-Rothman[1]
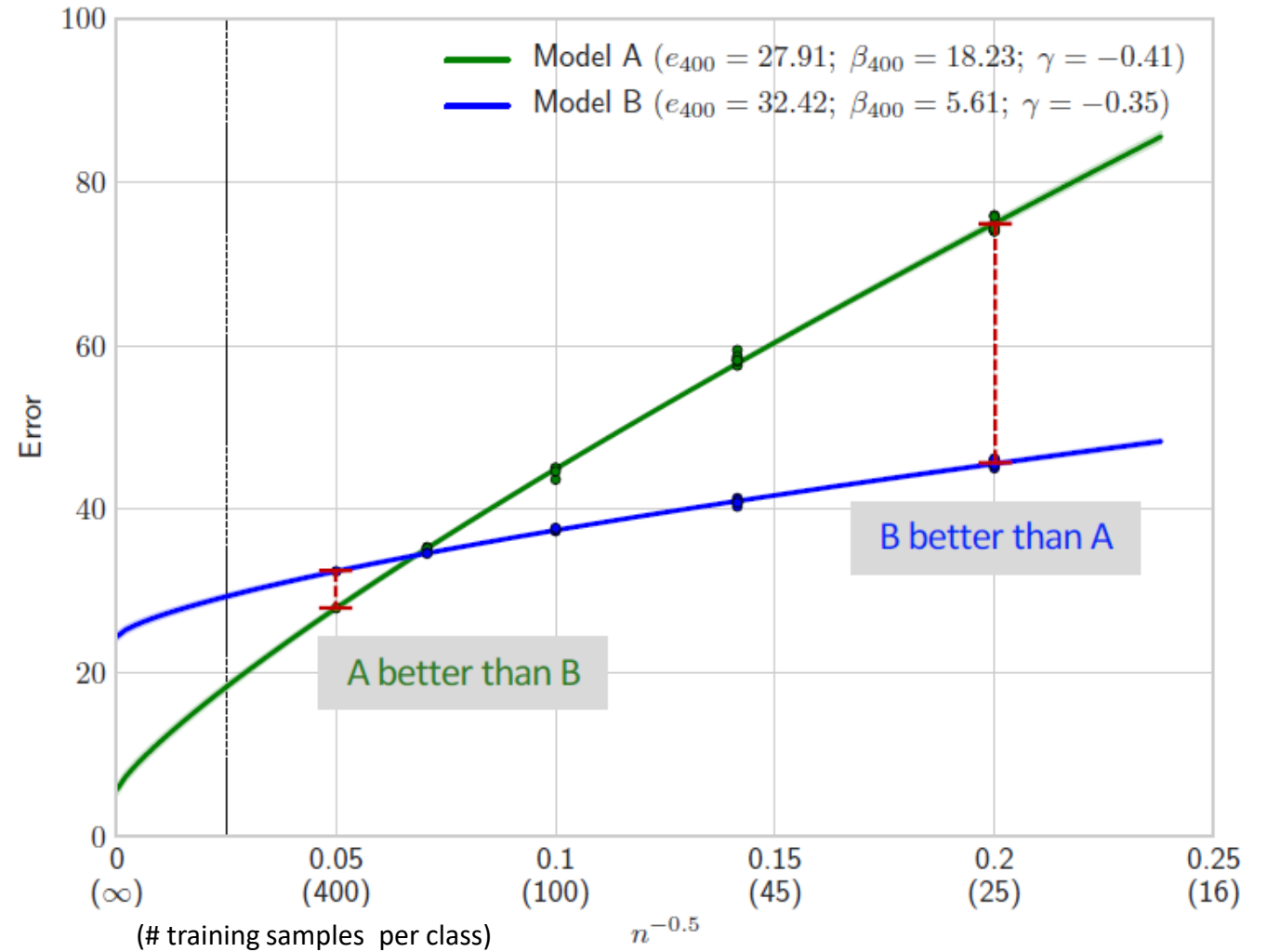
1

2

ILLINOIS
UNIVERSITY OF ILLINOIS AT URBANA-CHAMPAIGN

AI2
Allen Institute for AI

# Which classifier is better?

| Model | Error after full training (n=400) |
|-------|-----------------------------------|
| A     | 27.9%                             |
| B     | 32.4%                             |



Error

A better than B

0    0.05    0.1    0.15    0.2    0.25
($\infty$)  (400)  (100)  (45)   (25)   (16)

(# training samples  per class)          $n^{-0.5}$

# Which classifier is better?

| Model | Error after full training (n=400) |
|-------|-----------------------------------|
| A     | 27.9%                             |
| B     | 32.4%                             |



Model A ($e_{400} = 27.91$; $\beta_{400} = 18.23$; $\gamma = -0.41$)

Model B ($e_{400} = 32.42$; $\beta_{400} = 5.61$; $\gamma = -0.35$)

B better than A

A better than B

Error

$$n^{-0.5}$$

(# training samples per class)

0   0.05   0.1   0.15   0.2   0.25
($\infty$)  (400)  (100)  (45)  (25)  (16)

# Better characterize classifier performance with learning curve and measure of data reliance

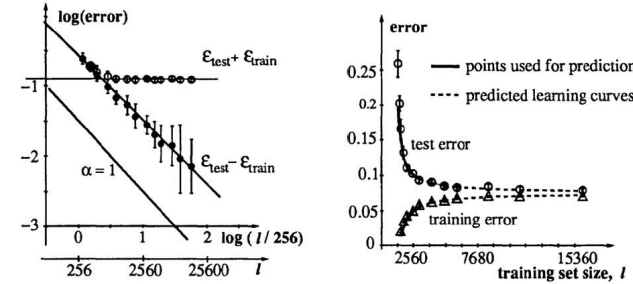| Model | Error after full training (n=400) | Data Reliance |
|-------|-----------------------------------|---------------|
| A     | 27.9%                             | 18.2          |
| B     | 32.4%                             | 5.6           |

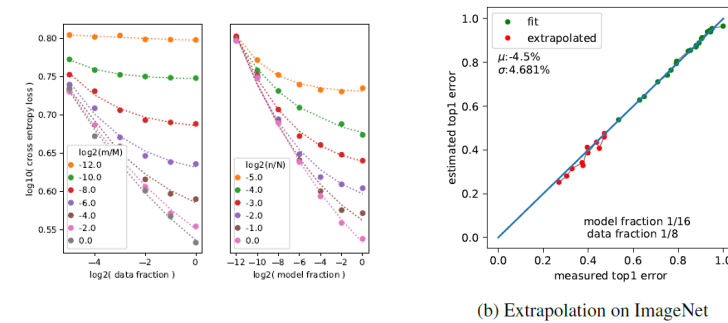# Learning curves have been shown useful, but there is no established methodology for how to use them in evaluation

**model selection**

Cortes et al.
NIPS 1993



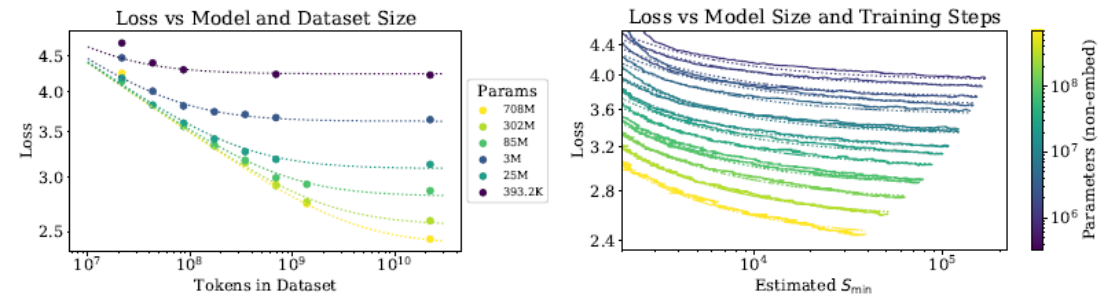**relationship analysis**
of model size, training size, computation

Rosenfeld et al.
ICLR 2020



**extrapolation**

Kaplan et al.
arxiv 2020

# Our goal: make it easy to improve classifier evaluations with learning curves

- Show how to model, fit, and display without using a lot of computation or paper space

- Show that learning curves provide useful insights

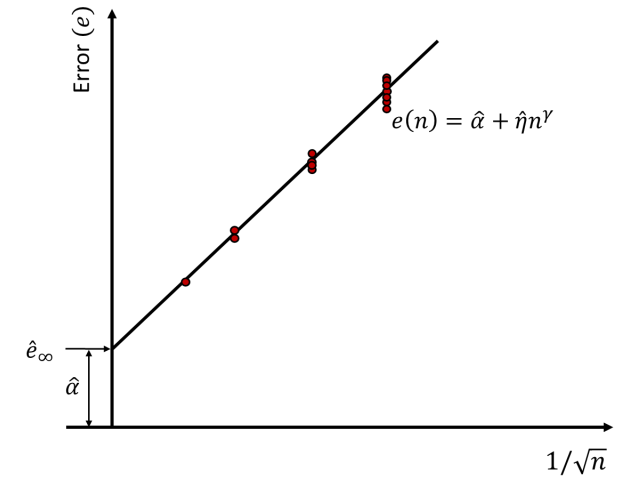# Model learning curves with extended power law

$$e(n) = \alpha + \eta n^{\gamma}$$

$n$: Number of training samples (per class)

$e$: Test error

## Well supported by

- Theory: bias-variance trade-off, many generalization bounds
- Practice: Hestness et al. 2017, Johnson & Nguyen 2017, Kaplan et al. 2020, Rosenfeld et al. 2020
- Our experiments

# Fit learning curves with weighted least squares


Error (e)
$e(n) = \hat{\alpha} + \hat{\eta} n^\gamma$
$\hat{e}_\infty$
$\hat{\alpha}$
$1/\sqrt{n}$

1. Given $\gamma$, solve for $\alpha, \eta$

$$\mathcal{G}(\gamma) = \min_{\alpha, \eta} \sum_{i=1}^{S} \sum_{j=1}^{F_i} w_{ij} \left(e_{ij} - \alpha - \eta n^\gamma\right)^2$$

$e_{ij}$: observation test error on split $i$ with size $F_i$

$w_{ij}$: accounts for variance of $e_{ij}$ and number of splits
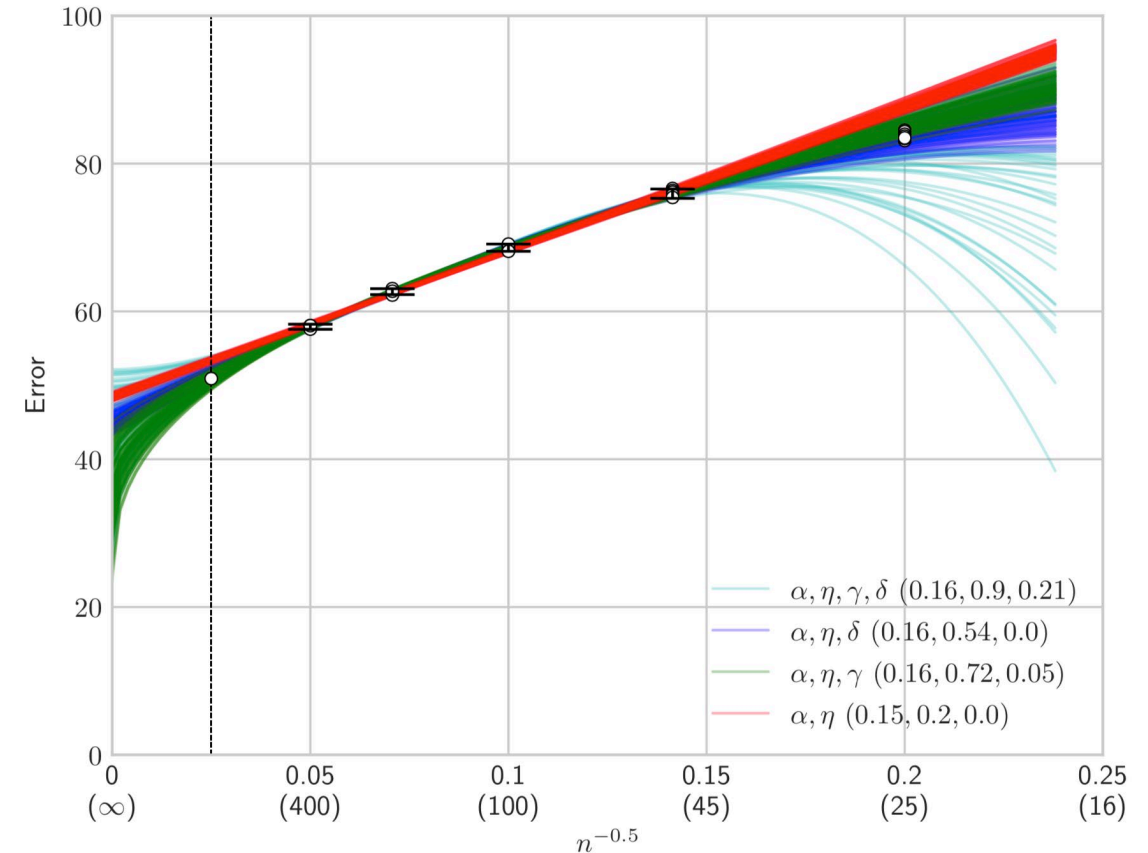
2. Step over $\gamma$

$$\min_{\gamma \in (-1, 0)} \mathcal{G}(\gamma) + \lambda |\gamma + 0.5|$$

# Extended power law model and weighted fitting lead to better prediction of error and more stable parameters

| Functional Form | Parameters |
|---|---|
| $e(n) = \alpha + \eta n^{-0.5}$ | $\alpha, \eta$ |
| $e(n) = \alpha + \eta n^{\gamma}$ | $\alpha, \eta, \gamma$ |
| $e(n) = \alpha + \eta n^{-0.5} + \delta n^{-1}$ | $\alpha, \eta, \delta$ |
| $e(n) = \alpha + \eta n^{\gamma} + \delta n^{2\gamma}$ | $\alpha, \eta, \gamma, \delta$ |

| Params | Weights | $R^2$ | RMSE 25 | 50 | 100 | 200 | 400 | avg | p-value |
|---|---|---|---|---|---|---|---|---|---|
| $\alpha, \eta, \gamma$ | $\frac{1}{\sigma_i^2 F_i}$ | 0.998 | 2.40 | 0.86 | **0.54** | 0.57 | **0.85** | **1.04** | - |
| | $\frac{1}{\sigma_i^2}$ | 0.999 | **2.38** | 0.83 | 0.69 | 0.54 | 1.08 | 1.10 | 0.06 |
| | 1 | 0.998 | 2.66 | 0.86 | 0.79 | **0.50** | 1.26 | 1.21 | 0.008 |
| $\alpha, \eta$ | $\frac{1}{\sigma_i^2 F_i}$ | 0.988 | 3.41 | 1.09 | 0.69 | 0.72 | 1.21 | 1.42 | <0.001 |
| $\alpha, \eta, \delta$ | $\frac{1}{\sigma_i^2 F_i}$ | 0.999 | 2.89 | **0.74** | 0.68 | 0.56 | 0.94 | 1.16 | 0.05 |
| $\alpha, \eta, \delta, \gamma$ | $\frac{1}{\sigma_i^2 F_i}$ | 0.999 | 3.46 | **0.74** | 0.70 | 0.59 | 1.00 | 1.30 | 0.02 |



Legend:
$\alpha, \eta, \gamma, \delta$ (0.16, 0.9, 0.21)
$\alpha, \eta, \delta$ (0.16, 0.54, 0.0)
$\alpha, \eta, \gamma$ (0.16, 0.72, 0.05)
$\alpha, \eta$ (0.15, 0.2, 0.0)

Leave-one-train-size-out error analysis:
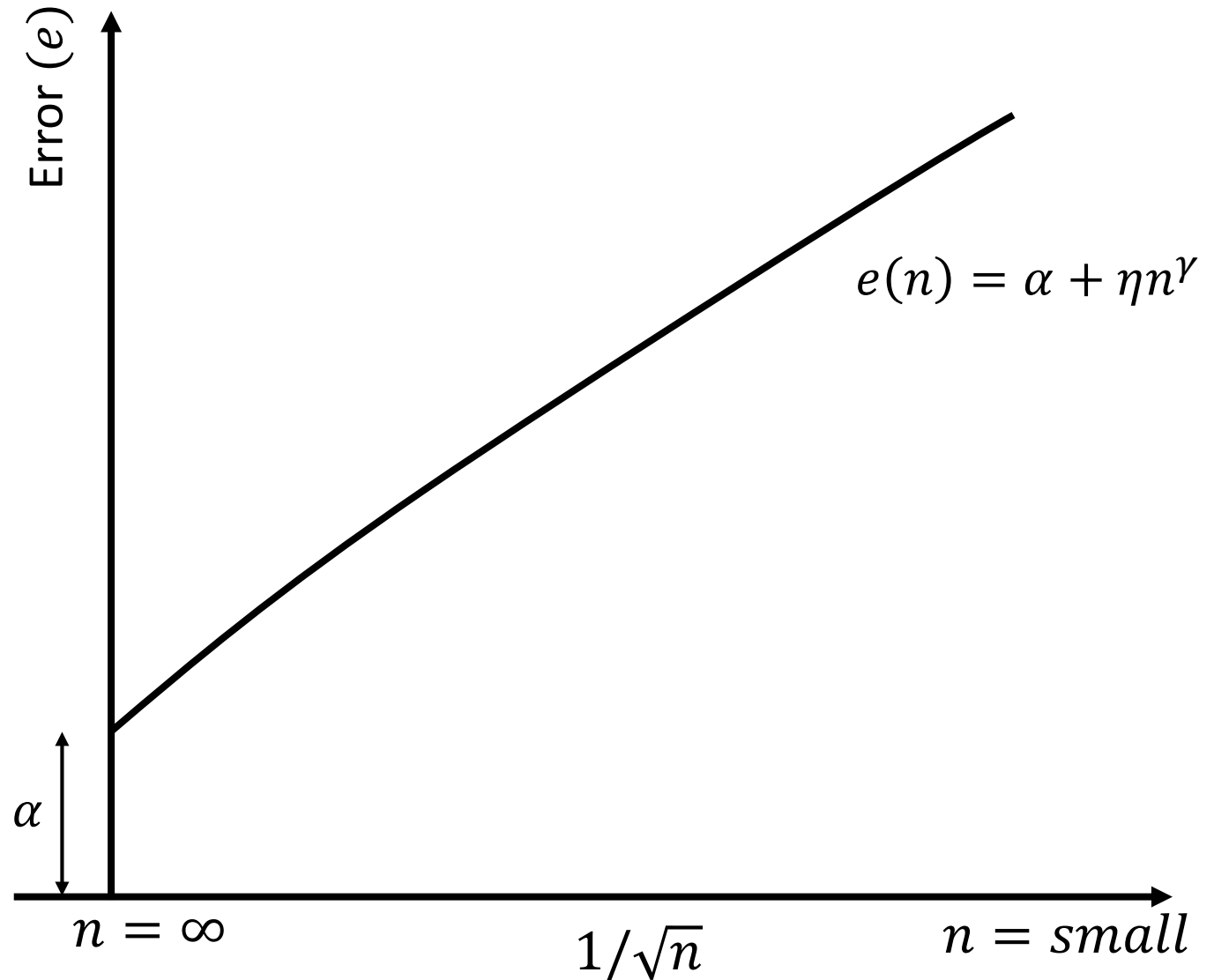Model accounts for 99.8% of *e(n)* variance and typically predicts on held out training size within 1% error

Stability analysis:
Given only 4 error observations (evaluations for 4 training sizes), our model better extrapolates and leads to more stable parameters with resampling

# Display learning curves as error vs $n^{-0.5}$

- $\gamma \approx -0.5$ typically

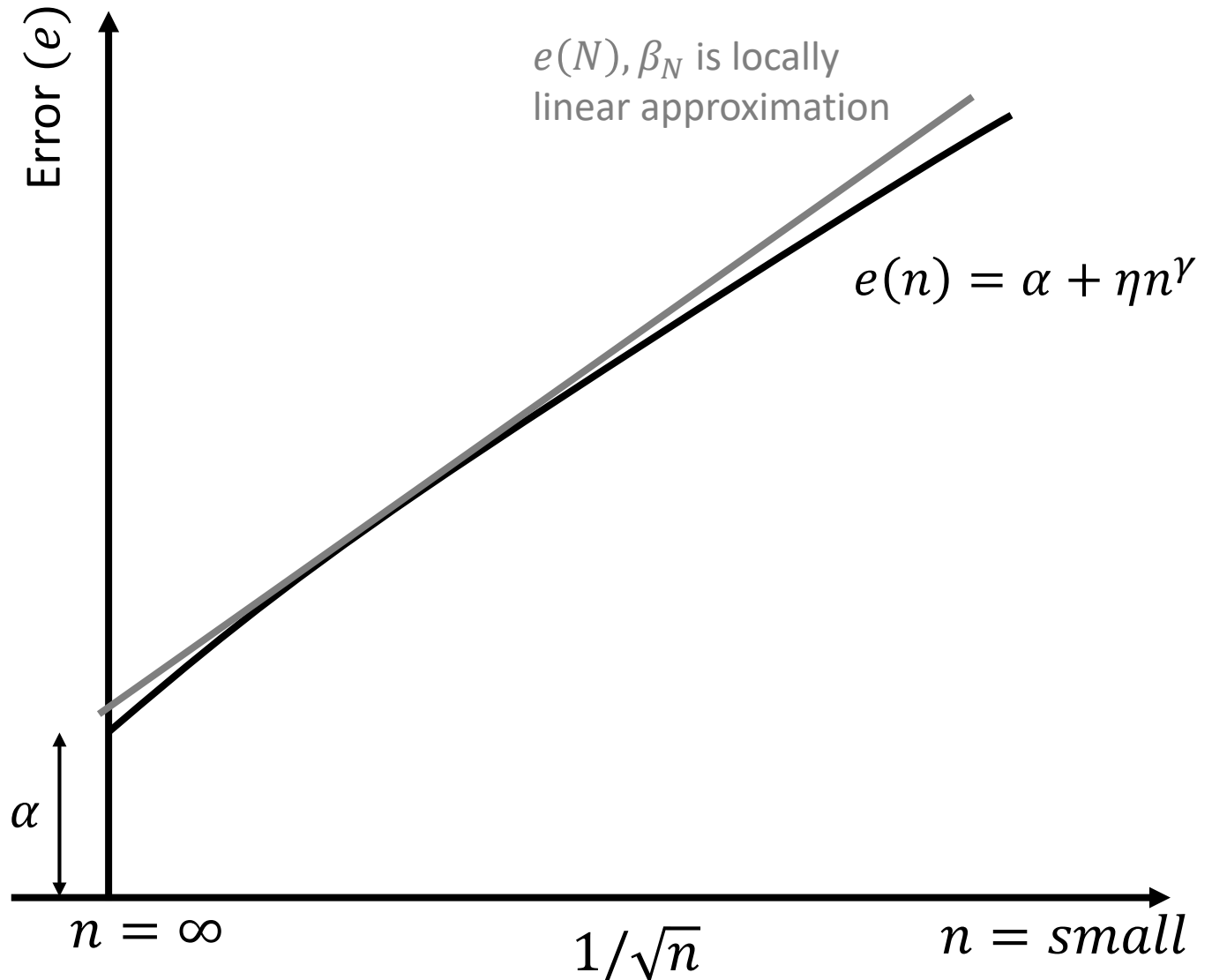- Curves are roughly linear in $n^{-0.5}$

- Can see full range of $n$



$$e(n) = \alpha + \eta n^\gamma$$

Error ($e$)

$\alpha$

$n = \infty$

$1/\sqrt{n}$

$n = small$

# How to characterize/summarize learning curves

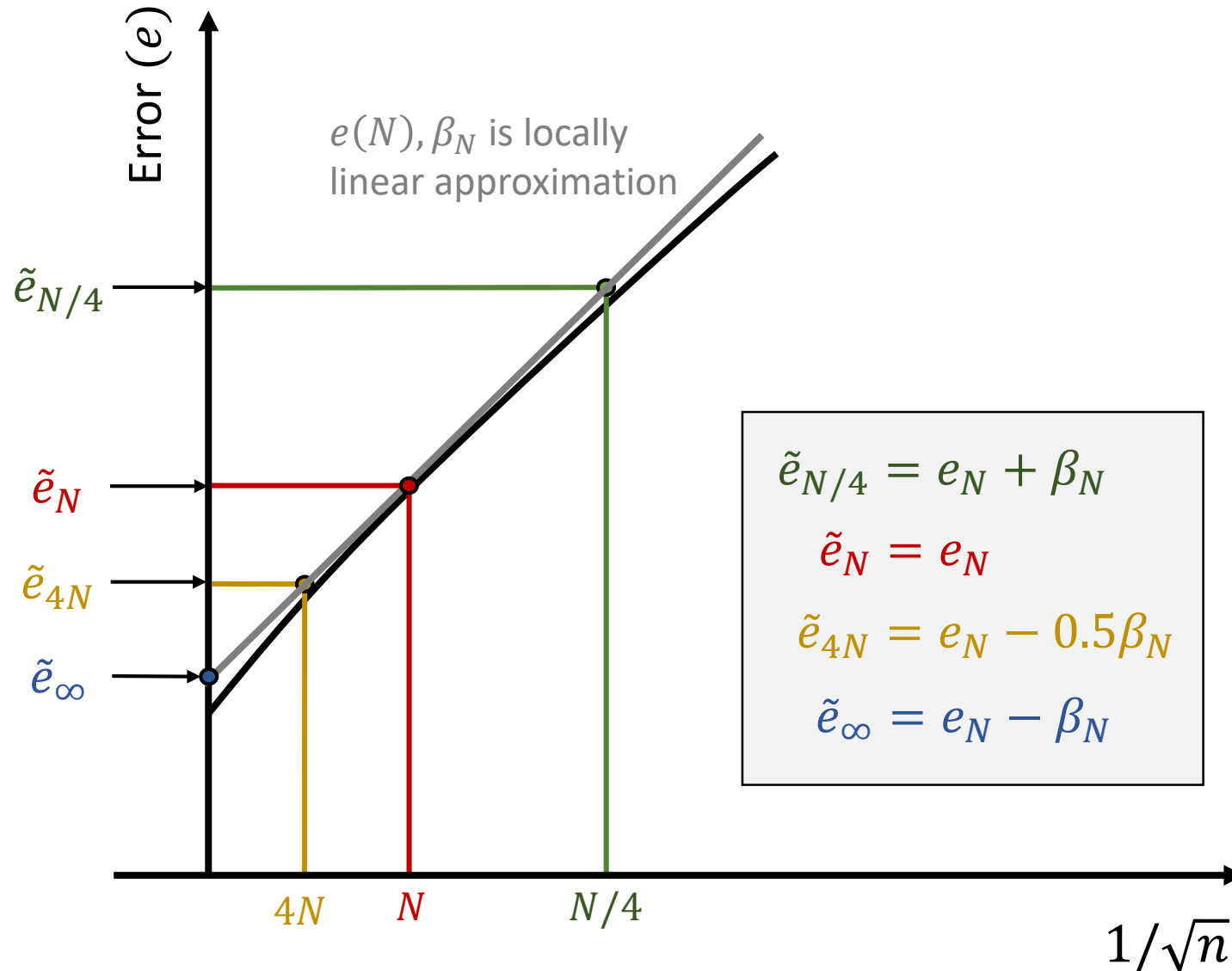**Problem:** $\gamma, \eta, \alpha$ highly covariant with observation perturbations and not individually comparable across curves

## Solution: re-parameterize

- $e_N = e(N)$ is error at full training size $N$
- $\beta_N = e'(N)/\sqrt{N}$ is data-reliance
- Can recover $\alpha, \gamma, \eta$ from $e_N, \beta_N, \gamma$

$e(N), \beta_N$ is locally linear approximation

$$e(n) = \alpha + \eta n^\gamma$$

Error $(e)$

$\alpha$

$n = \infty$

$1/\sqrt{n}$

$n = small$

# $\beta_N$ characterizes how error depends on data size, is stably estimated under perturbation, and easy to derive for other models
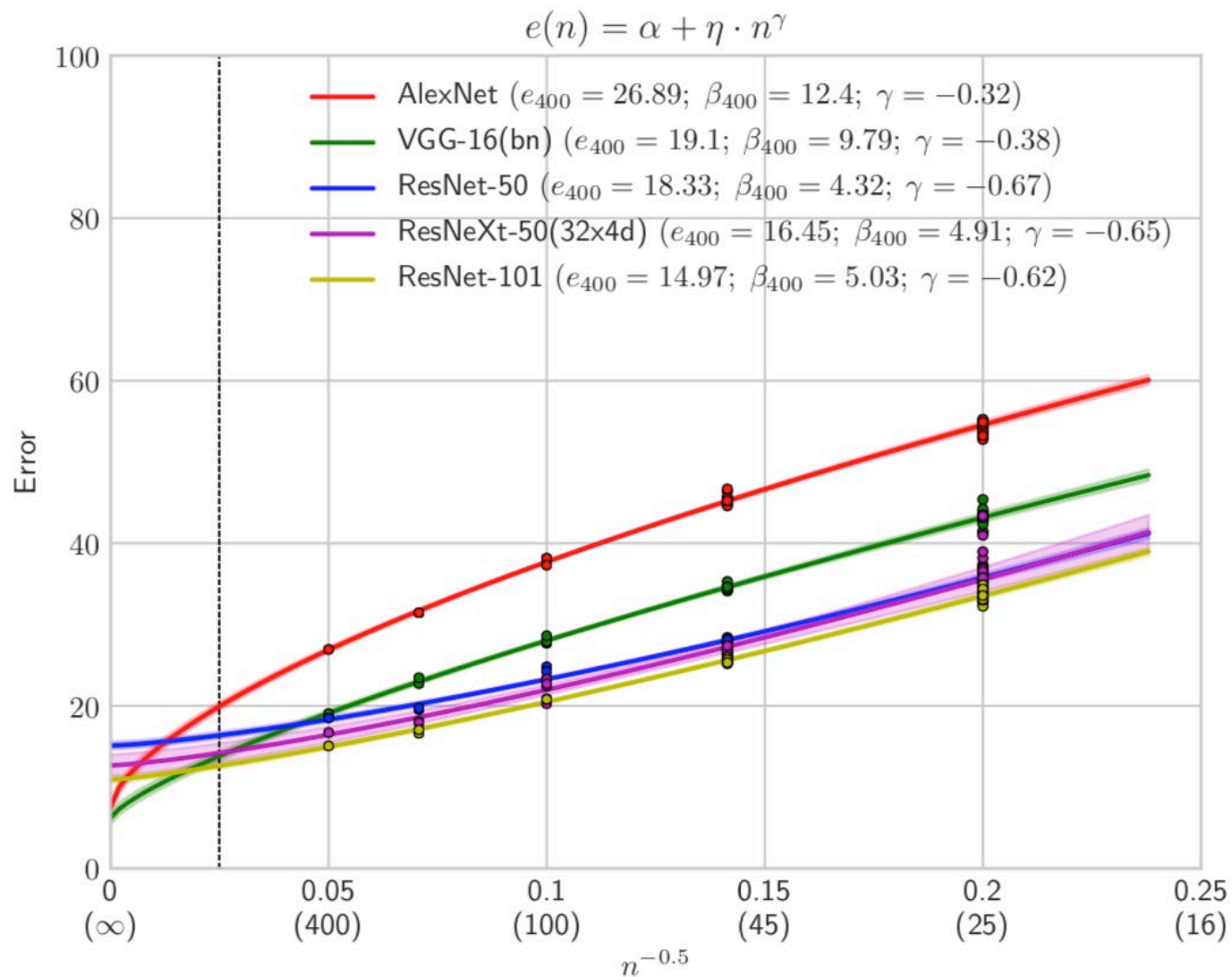


$e(N), \beta_N$ is locally linear approximation

$\tilde{e}_{N/4}$

$\tilde{e}_N$

$\tilde{e}_{4N}$

$\tilde{e}_\infty$

$$\tilde{e}_{N/4} = e_N + \beta_N$$

$$\tilde{e}_N = e_N$$

$$\tilde{e}_{4N} = e_N - 0.5\beta_N$$

$$\tilde{e}_\infty = e_N - \beta_N$$

Display learning curve analysis with only one extra column for $\beta_N$

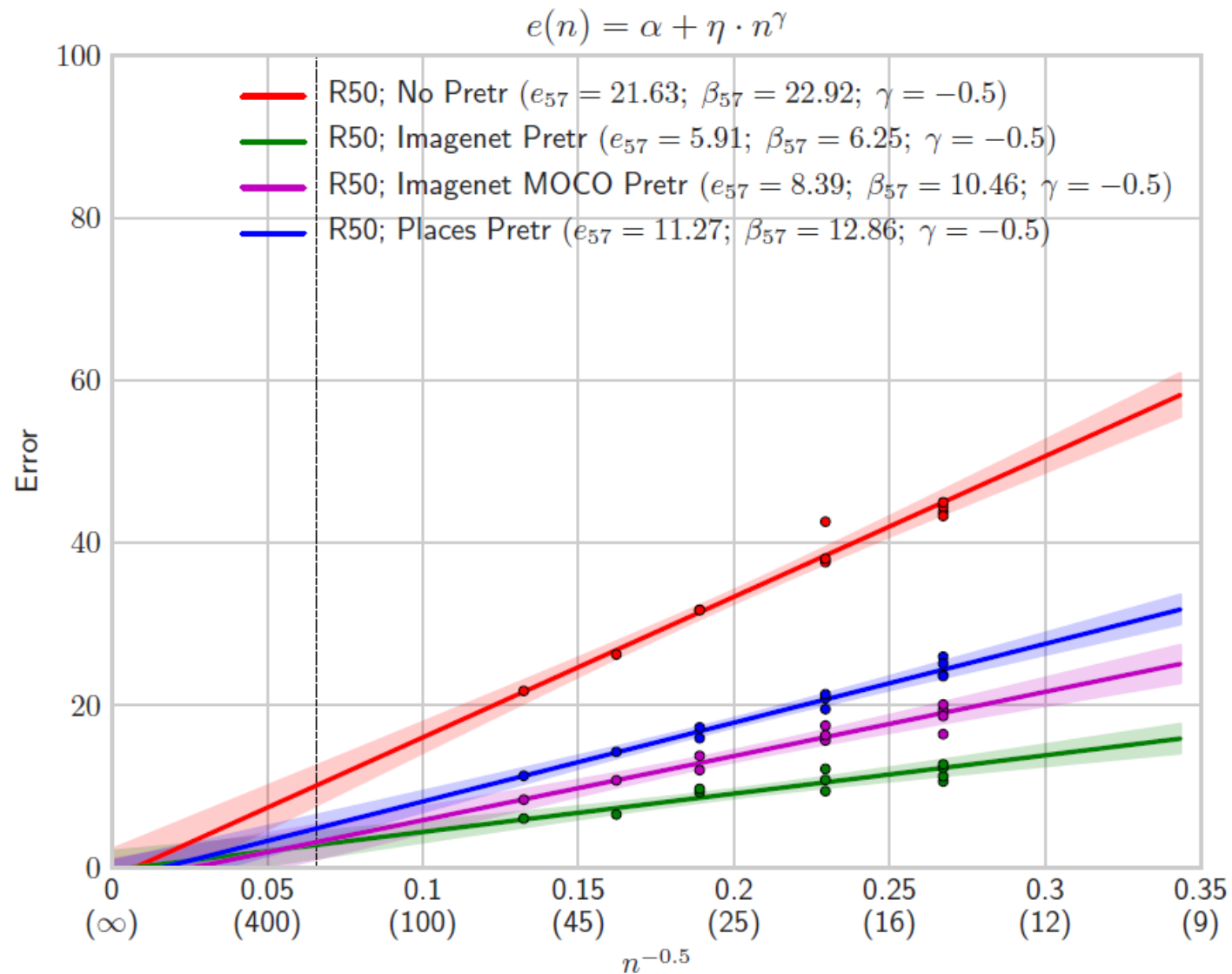| Model | $e_{400}$ | $\beta_{400}$ |
|-------|-----------|---------------|
| A | 27.91 | 18.23 |
| B | 32.42 | 5.61 |

**Comparing Architectures**
(Cifar-100, no pretraining)

More recent architectures
achieve lower error with
less data-reliance



$$e(n) = \alpha + \eta \cdot n^{\gamma}$$

— AlexNet $(e_{400} = 26.89;\ \beta_{400} = 12.4;\ \gamma = -0.32)$
— VGG-16(bn) $(e_{400} = 19.1;\ \beta_{400} = 9.79;\ \gamma = -0.38)$
— ResNet-50 $(e_{400} = 18.33;\ \beta_{400} = 4.32;\ \gamma = -0.67)$
— ResNeXt-50(32x4d) $(e_{400} = 16.45;\ \beta_{400} = 4.91;\ \gamma = -0.65)$
— ResNet-101 $(e_{400} = 14.97;\ \beta_{400} = 5.03;\ \gamma = -0.62)$

**Effect of Pretraining Source**
(Caltech-101)

Supervised ImageNet pretraining
provides much lower data-reliance
than others (for classification)



$$e(n) = \alpha + \eta \cdot n^\gamma$$

R50; No Pretr ($e_{57} = 21.63$; $\beta_{57} = 22.92$; $\gamma = -0.5$)
R50; Imagenet Pretr ($e_{57} = 5.91$; $\beta_{57} = 6.25$; $\gamma = -0.5$)
R50; Imagenet MOCO Pretr ($e_{57} = 8.39$; $\beta_{57} = 10.46$; $\gamma = -0.5$)
R50; Places Pretr ($e_{57} = 11.27$; $\beta_{57} = 12.86$; $\gamma = -0.5$)

**Effect of Depth**
(Cifar-100, ImageNet pretrained,
not fine-tuned)

Features from deeper networks
are better across range of
training sizes (deeper is almost
always better)



$$e(n) = \alpha + \eta \cdot n^{\gamma}$$

Resnet-18 ($e_{400} = 32.42$; $\beta_{400} = 5.61$; $\gamma = -0.35$)
Resnet-34 ($e_{400} = 30.02$; $\beta_{400} = 5.44$; $\gamma = -0.33$)
Resnet-50 ($e_{400} = 28.63$; $\beta_{400} = 6.66$; $\gamma = -0.22$)
Resnet-101 ($e_{400} = 25.86$; $\beta_{400} = 6.17$; $\gamma = -0.21$)

Error

$n^{-0.5}$

# Check out our deep learning quiz and additional experiments in the paper to test your beliefs

| **Popular beliefs** | Your guess | Supp-orted? | Exp. figures |
|---|---|---|---|
| *Pre-training* on similar domains nearly always helps compared to training from scratch. | ☐ | Y | 5a, 5b, 6 |
| *Pre-training*, even on similar domains, introduces bias that would harm performance with a large enough training set. | ☐ | U | 6 |
| *Self-/un-supervised training* performs better than supervised pre-training for small datasets. | ☐ | N | 6 |
| *Fine-tuning* the entire network (vs. just the classification layer) is only helpful if the training set is large. | ☐ | N | 5a, 5b |
| *Increasing network depth*, when fine-tuning, harms performance for small training sets, due to an overly complex model. | ☐ | N | 7a |
| *Increasing network depth*, when fine-tuning, is more helpful for larger training sets than smaller ones. | ☐ | N | 7a |
| *Increasing network depth*, if the backbone is frozen, is more helpful for smaller training sets than larger ones. | ☐ | N | 7d |
| *Increasing depth or width* improves more than ensembles of smaller networks with the same number of parameters. | ☐ | Y | 7f |
| *Data augmentation* is roughly equivalent to using a $m$-times larger training set for some $m$. | ☐ | Y | 8 |

# Use learning curves to better evaluate your research contributions

```python
from lc.measurements import CurveMeasurements
from lc.curve import LearningCurveEstimator
from omegaconf import OmegaConf
import matplotlib
import matplotlib.pyplot as plt

# Load error measurements
curvems = CurveMeasurements()
curvems.load_from_json('data/no_pretr_ft.json')

# Load config
cfg = OmegaConf.load('lc/config.yaml')

# Estimate curve
curve_estimator = LearningCurveEstimator(cfg)
curve, objective = curve_estimator.estimate(curvems)

# Plot
curve_estimator.plot(curve,curvems,label='No Pretr; Ft')
plt.show()
```
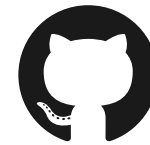
# Thank You

https://github.com/allenai/learning-curve

https://prior.allenai.org/projects/lcurve