# Continuous Integration of Machine Learning Models with `ease.ml/ci`
## *Towards a Rigorous Yet Practical Treatment*

**Cedric Renggli** (*PhD Candidate at DS3Lab, ETH Zurich*) - cedric.renggli@inf.ethz.ch

"Continuous Integration of Machine Learning Models: A Rigorous Yet Practical Treatment, Renggli et. al." at SysML 2019

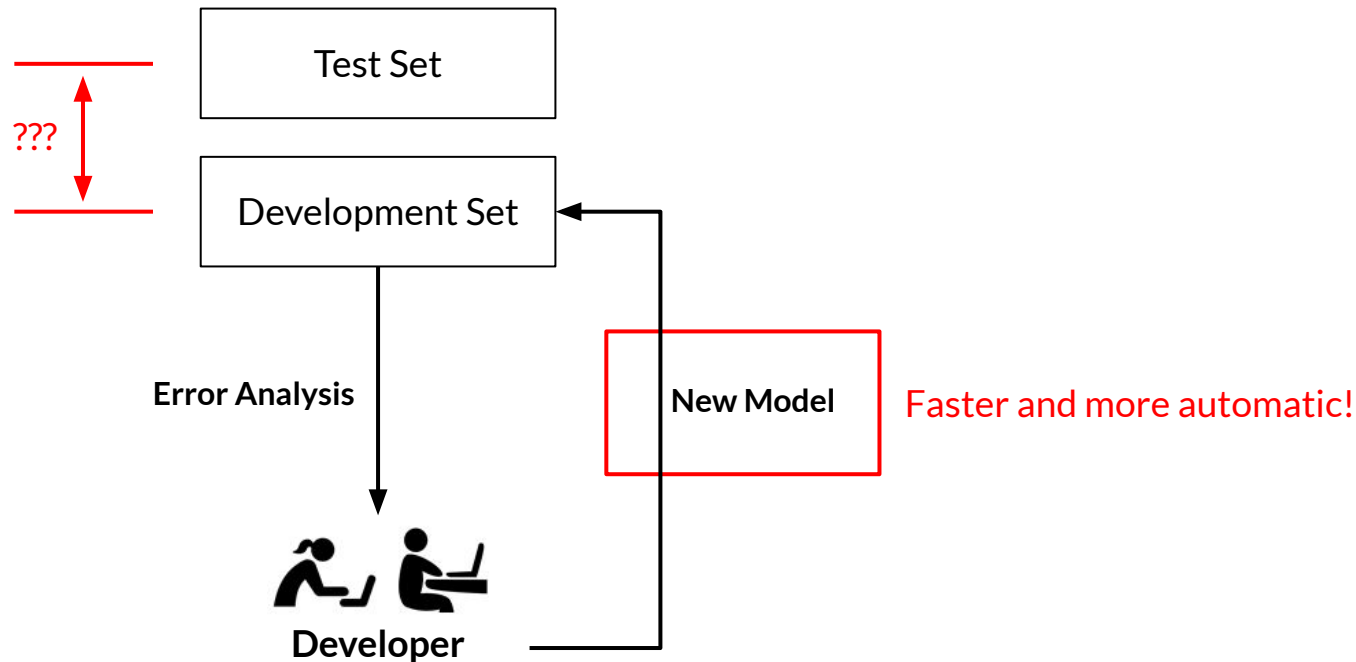# ML → Engineering Task

> **Andrew Ng** ✔
> @AndrewYNg
>
> 1/The rise of Software Engineering required inventing processes like version control, code review, agile, to help teams work effectively. The rise of AI & Machine Learning Engineering is now requiring new processes, like how we split train/dev/test, model zoos, etc.
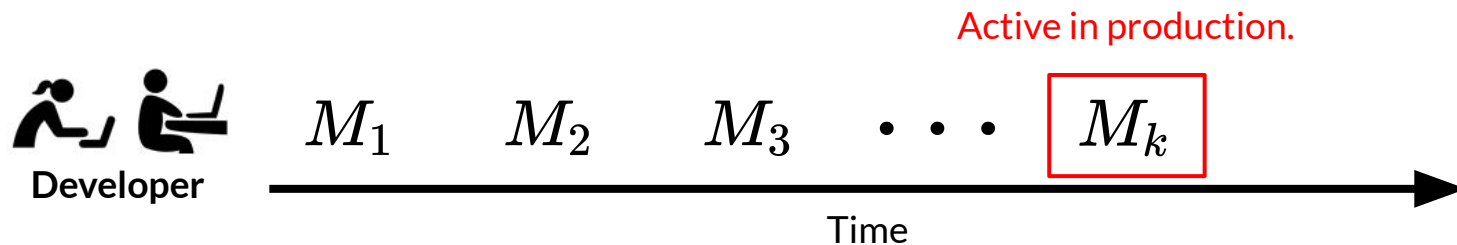>
> 6:59 PM · Jan 3, 2019 · Twitter Web Client
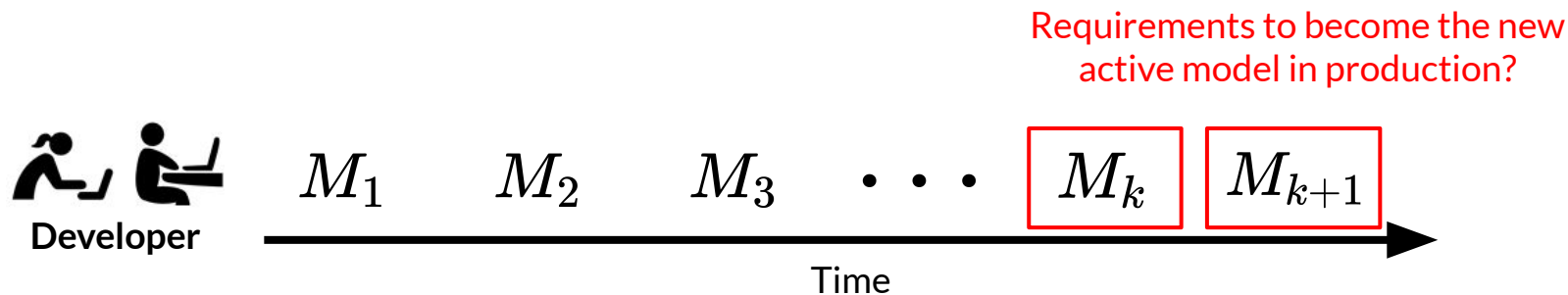>
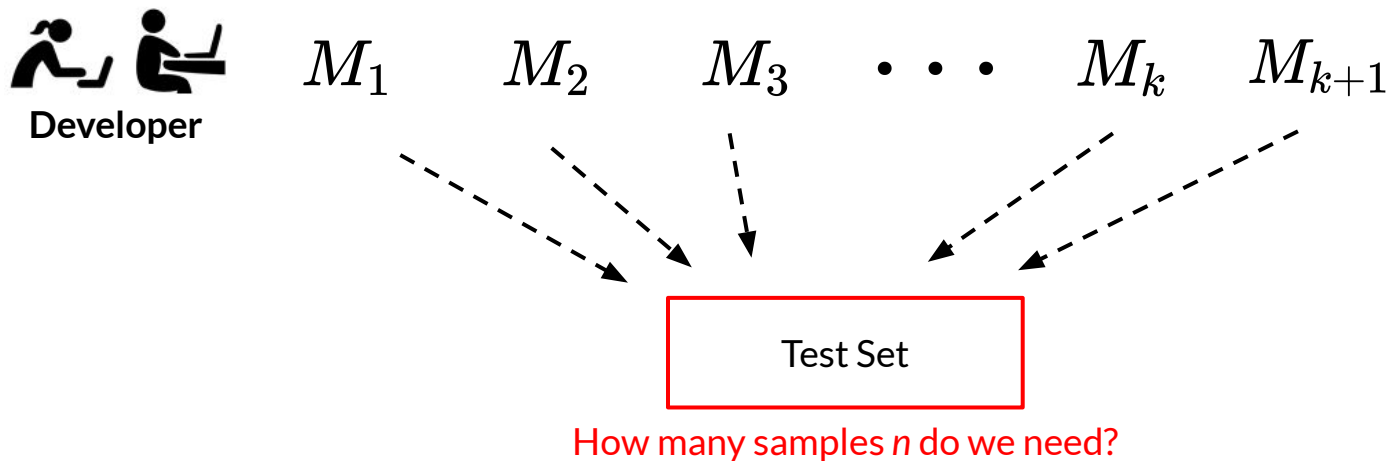> **1.1K** Retweets    **3.4K** Likes

# Typical ML Dev Process

Test Set

??? 

Development Set

Error Analysis

New Model

Faster and more automatic!

Developer

# Continuous Development of ML Models



Active in production.

$M_1 \quad M_2 \quad M_3 \quad \cdots \quad M_k$

Developer

Time

# Continuous Development of ML Models

Requirements to become the new active model in production?

$$M_1 \quad M_2 \quad M_3 \quad \cdots \quad M_k \quad M_{k+1}$$

**Developer**

Time

# Continuous Development of ML Models

$$M_1 \quad M_2 \quad M_3 \quad \cdots \quad M_k \quad M_{k+1}$$

**Developer**

Test Set

How many samples *n* do we need?

# System Overview



(5) When test labels lose statistical power, downgrade to val set and let developers know

(4) Ask for *n* test labels when it needs more

**ML Repo (e.g., Github)**

Encryption - Protected

**Manager**

Public

**Developer**

Public

(1) Specify Requirements

e.g., all models checked in should have accuracy > 0.8 $(\epsilon, \delta)$-approximation.

(2) Commit a stream of *T* models

(3) Receive Pass/Fail signal per commit

EASE.ML PASSING    EASE.ML FAILING

# Managers Specify Requirements

**Manager**

R1: New model needs to be better than the old model by at least 1%, with probability 0.999.

```
n - o > 0.01, p > 0.999
```

R2: New model cannot be different from the old model on more than 10% of predictions, with probability 0.999.

```
d < 0.1, p > 0.999
```

R3: New model always have accuracy higher than 0.8, with probability 0.999.

```
n > 0.8, p > 0.999
```

R4: Satisfy both R1 and R2, with probability 0.999.

```
n - o > 0.01 and d < 0.1, p > 0.999
```

# Developers Task



**Developer**

**Develop a ML model and <span style="color:red">commit</span>.**

# Developers Task



**Developer**

**Develop a new ML model and recommit.**

Core Technical Component:

*Adaptive Statistical Queries*

We are inspired by the following seminal work:

- The ladder: A reliable leaderboard for machine learning competitions. Blum and Hardt, 2015
- The algorithmic foundations of differential privacy. Dwork et. al., 2014
- The reusable holdout: Preserving validity in adaptive data analysis. Dwork et. al., 2015
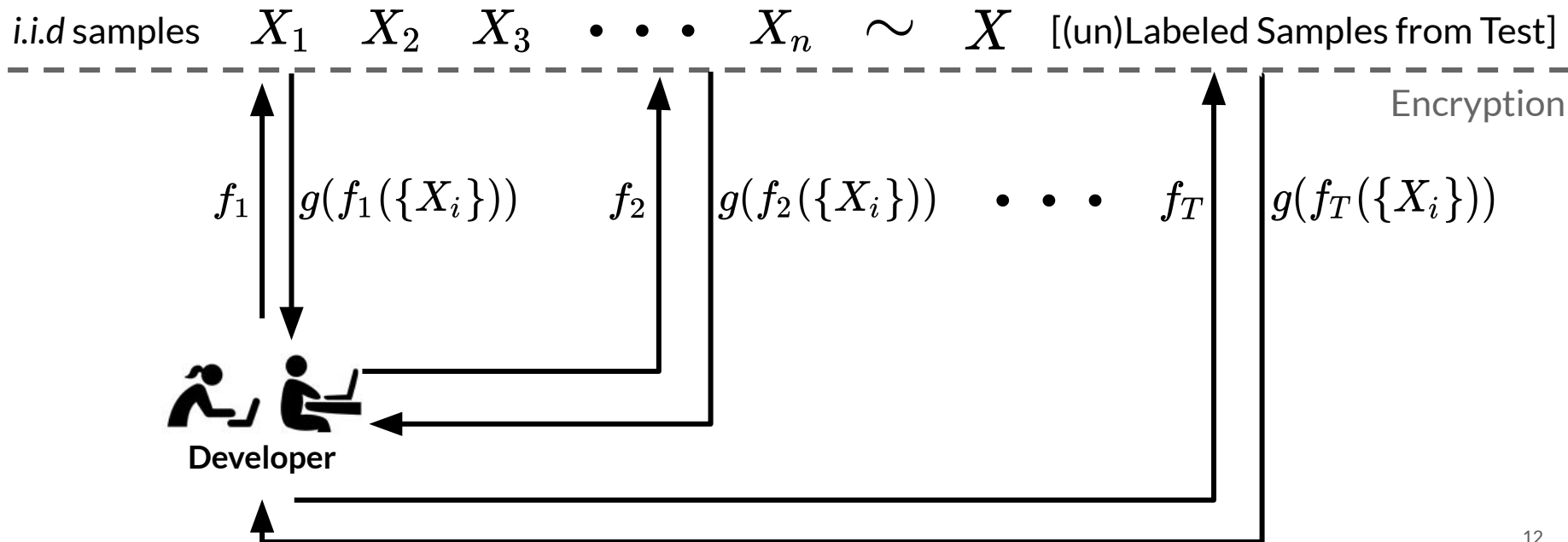
# Background: Adaptive Analytics

Contract between System and User:
$$\Pr\left[\exists t, |f_t(X_1,\ldots,X_n) - f_t(X)| > \epsilon\right] < \delta$$

Given $\varepsilon$, $\delta$, $T$, how large does $n$ need to be?

How can we decrease the dependency of $n$ on $\varepsilon$, $\delta$, $T$ as much as possible?

*i.i.d* samples   $X_1$   $X_2$   $X_3$   $\bullet\bullet\bullet$   $X_n$   $\sim$   $X$   [(un)Labeled Samples from Test]

Encryption

$f_1$   $g(f_1(\{X_i\}))$   $f_2$   $g(f_2(\{X_i\}))$   $\bullet\bullet\bullet$   $f_T$   $g(f_T(\{X_i\}))$

**Developer**

# Background: Single Steps – Hoeffding's Inequality

Theorem (Hoeffding, 1963):

Let $X_1, X_2, \ldots, X_n$ be i.i.d random variables with

$\forall X_i \; 0 \leq X_i \leq 1$ and $\overline{X} = \frac{1}{n} \sum_{i=1}^{n} X_i$ :

Then $\forall \epsilon$

$$\Pr\left[\overline{X} - \mathbb{E}[\overline{X}] \geq \epsilon\right] \leq \exp\left(-2n\epsilon^2\right).$$

$$\delta \leq \exp\left(-2n\epsilon^2\right) \quad \longrightarrow \quad n \geq \frac{\ln \frac{1}{\delta}}{2\epsilon^2}$$

# Background: Multiple Steps – Existing Solutions

$$f_2(\{X_i\}) = h_{g(f_1(\{X_1, X_2, \ldots, X_n\}))}(\{X_i\})$$

| Baseline Approach: Resampling | Ladder (Blum and Hardt, 2015) | Other DP - inspired approaches |
|---|---|---|
| *Require a new sample for each step.* | *Constrains how g(–) evolves over time.* | |

$\epsilon = 0.01$
$\delta = 0.001$
$T = 32$

$$n \geq T\frac{-\ln\frac{\delta}{T}}{2\epsilon^2} \approx 1.7M$$

$$n \geq 69K$$

**Expensive: ~53K / Day**

**g(-) is non-monotonic**

**Unclear how to add noise to g(-) in CI**

**Goal: Optimizing Sample Complexity for the *specific* regime that *our system cares about*.**

# Overview of Optimizations

> Goal: Optimizing Sample Complexity for the *specific* regime that *our system cares about*.

1) **General Optimization**

2) **Stable Signal**

3) **Conditional Variance**

4) **Active Labeling**

# Adaptive Analytics - Observation

| Observation: Not all labels are useful |
|:---:|

*Focus:* `n - o > 0.01, p > 0.999`



Old Model:    0    1    1    1    0

New Model:    0    1    1    0    1

Same predictions – Not useful
to estimate the difference

*If new models and old models are only different in their prediction with probability $v$, how many savings can we have in terms of labels (NOT SAMPLES) that we need to provide?*

If the probability of two models being different is $v \sim O(\sqrt{\varepsilon})$, than the amount of labels we need is $n \geq O(1/\varepsilon)$.

Hoeffding          **15K** samples/signal
$v$ = 0.1            **2.2K** samples/signal
(Assuming unlabeled data points are free)

# `ease.ml/ci` in Action

## ease.ml/ci

```
$ git commit -m newmodel
```



## # of Labels/32 Models

**Popular Use Cases:** ($\varepsilon = 0.0125$)

| | **Baseline** | **`ease.ml/ci`** |
|---|---|---|
| `n - o > 0.01 and d < 0.1` | **4.8M** (150K / Day) | **41K** (1.3K / Day) |
| `n > 0.8` | **1.1M** (35K / Day) | **95K** (3K / Day) |

**Cheap Mode:** ($\varepsilon = 0.025$)

| | | |
|---|---|---|
| `n - o > 0.01 and d < 0.1` | **1.2M** (38K / Day) | **11K** (330 / Day) |
| `n > 0.8` | **283K** (8.9K / Day) | **24K** (745 / Day) |

**10s / Label**

300 Labels / Day => < 1 Hour / Day

# `ease.ml/ci` in Action

## ease.ml/ci

```
$ git commit -m newmodel
```

EASE.ML  **EVALUATING**

EASE.ML  **FAILING**

EASE.ML  **PASSING**

**# of Labels/32 Models**

**Popular Use Cases:** ($\varepsilon = 0.0125$)

If ML is "Software 2.0", what are the other missing principles in "**Software Engineering 2.0**"?

`n - o > 0.01 and d < 0.1`

`n > 0.8`

| | Baseline | ease.ml/ci |
|---|---|---|
| | **4.8M** (150K / Day) | **41K** (1.3K / Day) |
| | **1.1M** (35K / Day) | **95K** (3K / Day) |
| | **1.2M** (38K / Day) | **11K** (330 / Day) |
| | **283K** (8.9K / Day) | **24K** (745 / Day) |

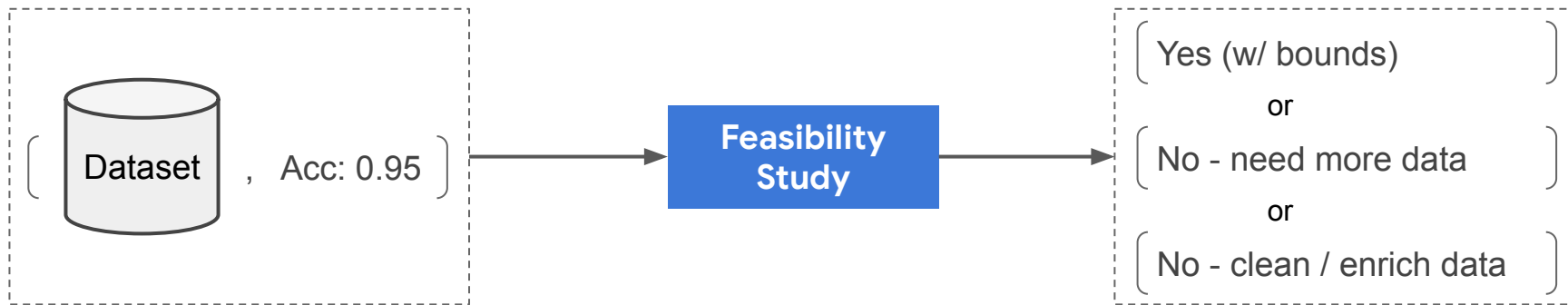**10s / Label**   300 Labels / Day => < 1 Hour / Day

18

# Teaser - Feasibility Study for ML Application

If ML is "Software 2.0", what are the other missing principles in "**Software Engineering 2.0**"?

Dataset , Acc: 0.95 → **Feasibility Study** → Yes (w/ bounds)

or

No - need more data

or

No - clean / enrich data

# Question?