




# Ukraine International Airlines Flight 752

[https://en.wikipedia.org/wiki/Ukraine\\_International\\_Airlines\\_Flight\\_752](https://en.wikipedia.org/wiki/Ukraine_International_Airlines_Flight_752)

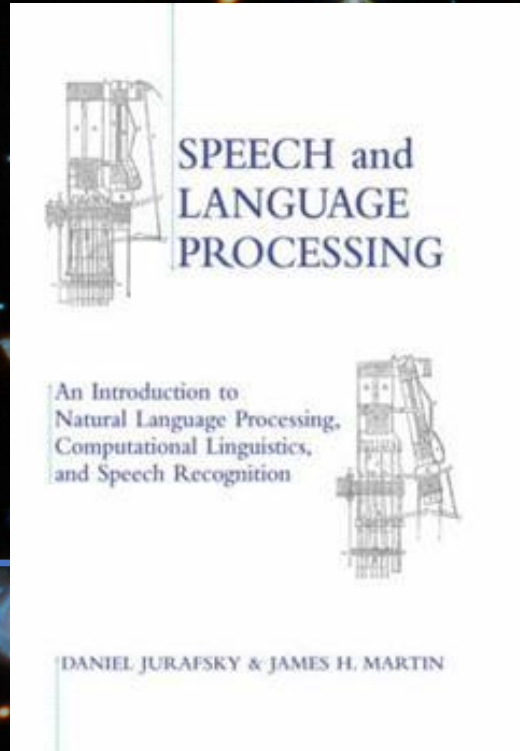


A medium shot of Keanu Reeves sitting on a blue couch. He is wearing a dark suit, a dark shirt, and a dark tie. He has long dark hair and a beard. The background is a dimly lit room with ornate silver-colored wall panels and a wooden bookshelf. A small microphone is clipped to his tie.

I Know that the ones who love us will miss us

[www.youtube.com/watch?v=etlBZInTE-I](https://www.youtube.com/watch?v=etlBZInTE-I)





# Naive Bayes and Sentiment Classification

---

## CH04



# Evaluate a Classifier

---





*"We're addicted and annotated data is our heroine."*

—?

---

Supervise

# Gold Labels

aka. Golden Truth, Golden Standard

---

Human-defined classes/labels for each document

Human judgment

Manually labeled

Manually annotated





Gold Labels

---

to err is human!





# Gold Labels

---

to err is human; to forgive, divine!

*"Save This Word! All people commit sins and make mistakes. God forgives them, and people are acting in a godlike (divine) way when they forgive."* - An Essay on Criticism, Alexander Pope.



# Silver Labels

aka. Silver Truth, Silver Standard

---

Gold is very expensive!

Finding gold is needs a lot of effort!

automated-defined classes/labels for each document

Machine judgment

Machine labeled

Machine annotated



# Transduction

Transductive Inference

---

Data has the labels already!

Language Models



# Evaluation

---

$$(x, y) \rightarrow f(x) == y$$

A deep space image showing a vast field of galaxies and stars against a black background. The galaxies are in various colors, including blue, orange, and red, and are scattered across the frame. Some are bright and clear, while others are faint and distant. The stars are small, bright points of light, some with visible diffraction patterns.

Evaluation

---

Boolean (Binary) Classifier



# Contingency Table

aka. Confusion Matrix



	Gold Positive	Gold Negative
<u>Model Positive</u>	True Positive	False Positive Type 1
Model Negative	False Negative Type 2	True Negative

# Perfect Classifier

---

	Gold Positive	Gold Negative
Model Positive	N+	0
Model Negative	0	N-



# Other Metrics

		True condition			
Total population		Condition positive	Condition negative	Prevalence = $\frac{\sum \text{Condition positive}}{\sum \text{Total population}}$	Accuracy (ACC) = $\frac{\sum \text{True positive} + \sum \text{True negative}}{\sum \text{Total population}}$
Predicted condition	Predicted condition positive	True positive	False positive, Type I error	Positive predictive value (PPV), Precision = $\frac{\sum \text{True positive}}{\sum \text{Predicted condition positive}}$	False discovery rate (FDR) = $\frac{\sum \text{False positive}}{\sum \text{Predicted condition positive}}$
	Predicted condition negative	False negative, Type II error	True negative	False omission rate (FOR) = $\frac{\sum \text{False negative}}{\sum \text{Predicted condition negative}}$	Negative predictive value (NPV) = $\frac{\sum \text{True negative}}{\sum \text{Predicted condition negative}}$
		True positive rate (TPR), Recall, Sensitivity, probability of detection, Power $= \frac{\sum \text{True positive}}{\sum \text{Condition positive}}$	False positive rate (FPR), Fall-out, probability of false alarm $= \frac{\sum \text{False positive}}{\sum \text{Condition negative}}$	Positive likelihood ratio (LR+) = $\frac{\text{TPR}}{\text{FPR}}$	Diagnostic odds ratio (DOR) = $\frac{\text{LR+}}{\text{LR-}}$  $F_1 \text{ score} = 2 \cdot \frac{\text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}}$
		False negative rate (FNR), Miss rate $= \frac{\sum \text{False negative}}{\sum \text{Condition positive}}$	Specificity (SPC), Selectivity, True negative rate (TNR) $= \frac{\sum \text{True negative}}{\sum \text{Condition negative}}$	Negative likelihood ratio (LR-) = $\frac{\text{FNR}}{\text{TNR}}$	

# Precision

High vs. Low

		gold standard labels		
		gold positive	gold negative	
system output labels	system positive	true positive	false positive	precision = $\frac{tp}{tp+fp}$
	system negative	false negative	true negative	
		recall = $\frac{tp}{tp+fn}$	spam	accuracy = $\frac{tp+tn}{tp+fp+tn+fn}$

*Handwritten notes: A red arrow points from 'system positive' to 'true positive'. Another red arrow points from 'false positive' to the precision formula. A red circle highlights the 'precision' formula. A red circle highlights the 'accuracy' formula. A red arrow points from 'recall' to the 'tp' in its numerator. A red arrow points from 'spam' to the 'fp' in the precision formula's denominator. A red arrow points from 'no + spam' to the 'tn' in the accuracy formula's numerator.*

**Figure 4.4** Contingency table

What scenarios require high precision?



# Recall

High vs. Low

		gold standard labels		
		gold positive	gold negative	
system output labels	system positive	true positive	false positive	precision = $\frac{tp}{tp+fp}$
	system negative	false negative	true negative	
		recall = $\frac{tp}{tp+fn}$		accuracy = $\frac{tp+tn}{tp+fp+tn+fn}$

*Handwritten notes:*

- Spam, not Spam* (above gold negative column)
- missile* (circled around true negative)
- release, not -* (next to precision formula)

**Figure 4.4** Contingency table

What scenarios require high recall?

# Precision-Recall





# Precision-Recall



# Balance Classes

~50% Positive, ~50% Negative

---

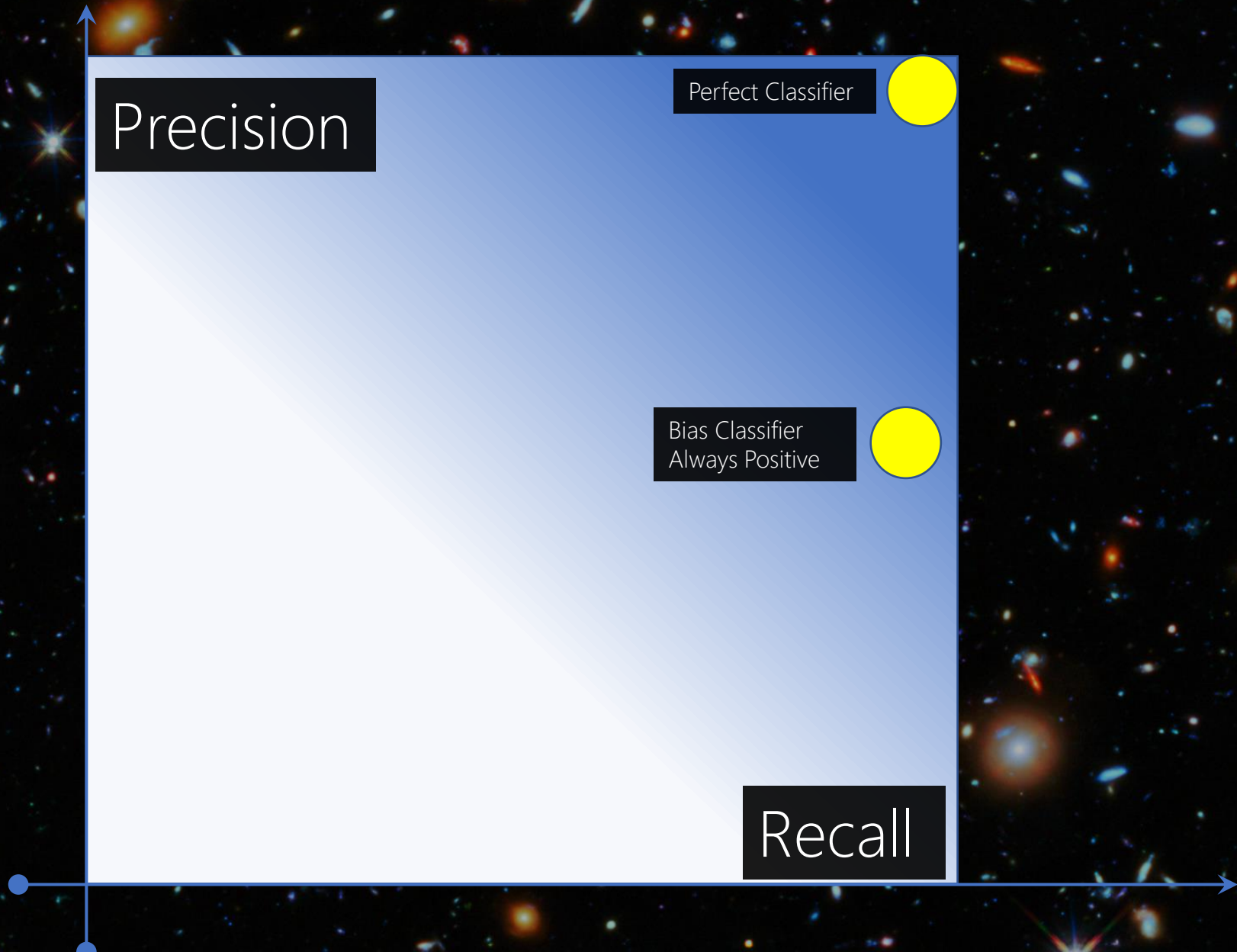
	Gold Positive (50)	Gold Negative (50)
Model Positive	50	50
Model Negative	0	0

$$\text{Precision} = \frac{50}{50+50} = 0.5$$

$$\text{Recall} = \frac{50}{50+0} = 1.0$$



# Precision-Recall



# Balance Classes

~50% Positive, ~50% Negative

---

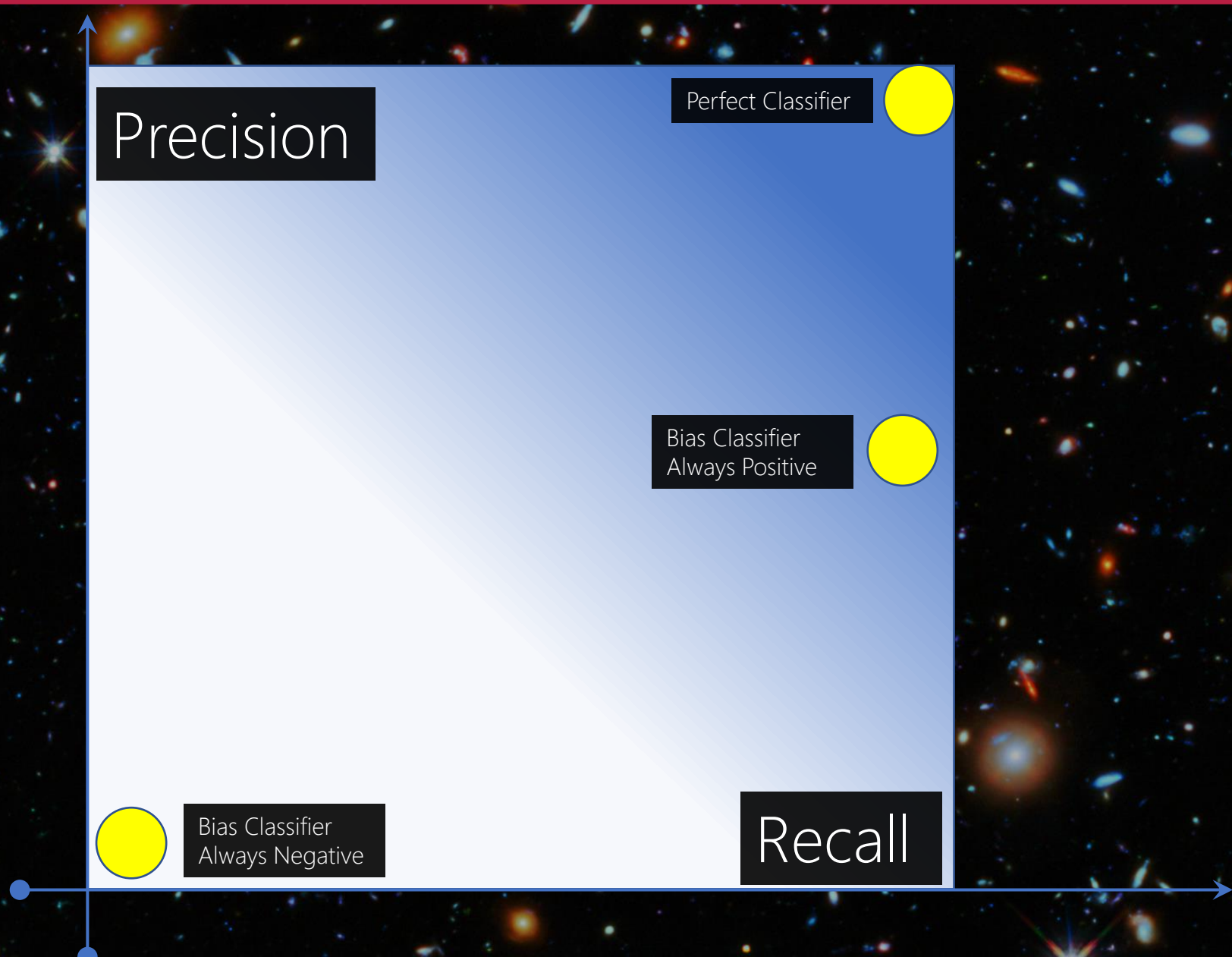
	Gold Positive (50)	Gold Negative (50)
Model Positive	0	0
Model Negative	50	50

$$\text{Precision} = \frac{0}{0+0} = 0.0$$

$$\text{Recall} = \frac{0}{50+0} = 0.0$$



# Precision-Recall



# Balance Classes

~50% Positive, ~50% Negative

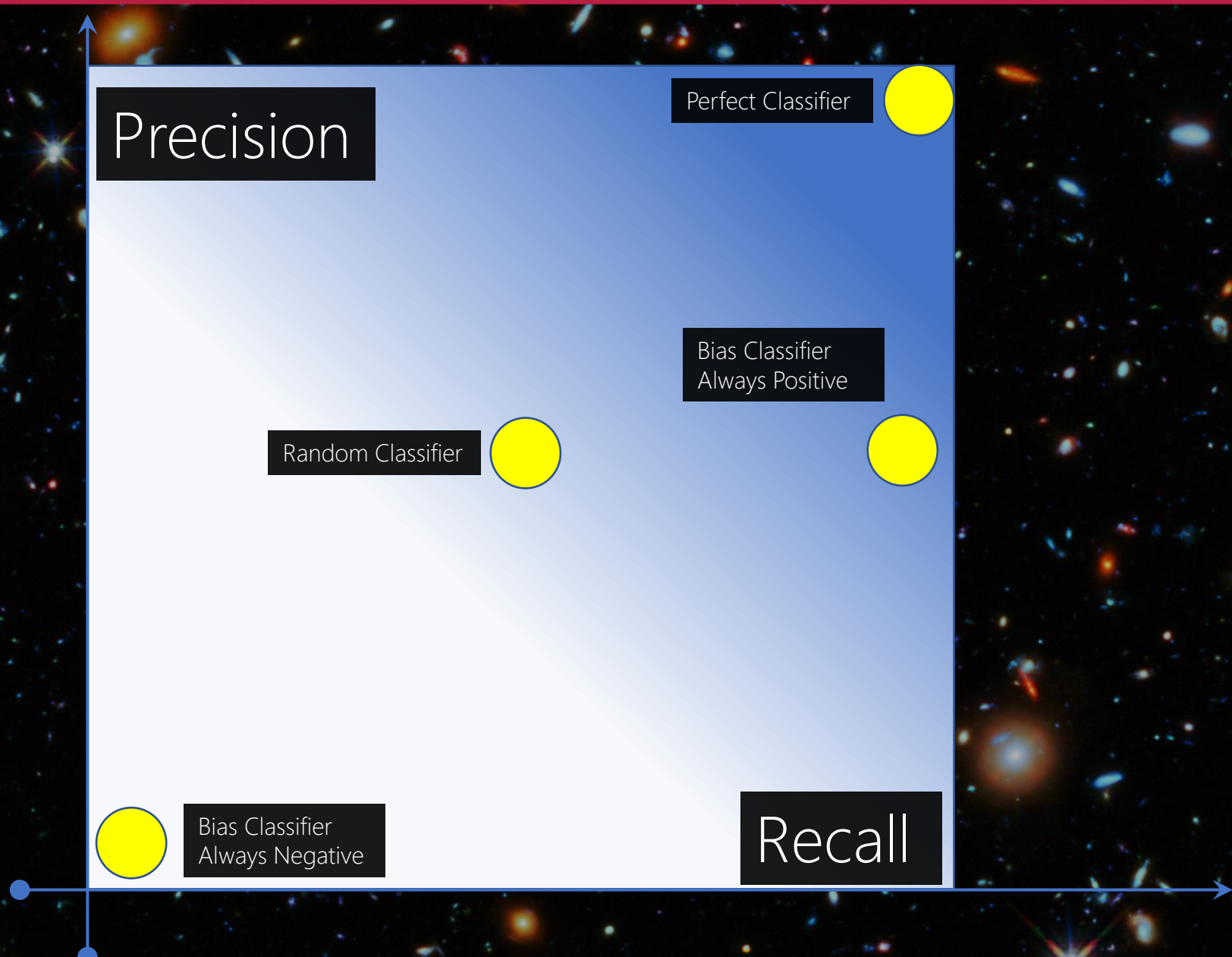
	Gold Positive (50)	Gold Negative (50)
Model Positive	25	25
Model Negative	25	25

$$\text{Precision} = \frac{25}{25+25} = 0.5$$

$$\text{Recall} = \frac{25}{25+25} = 0.5$$



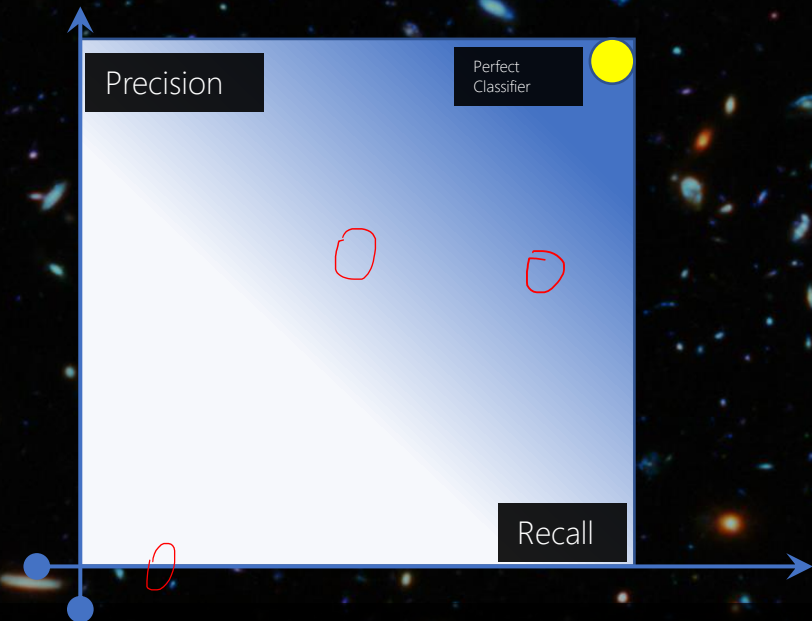
# Precision-Recall



# Imbalance (Unbalanced) Classes

~10% Positive, ~90% Negative

	Gold Positive (10)	Gold Negative (90)
Model Positive	?	?
Model Negative	?	?



Bias Positive Classifier?  
Bias Negative Classifier?  
Random Classifier?



# Average of Precision and Recall: A Single Metric

---

$$\text{AVG-PR} = \frac{P+R}{2} = 6.5$$

Same weights to Precision and Recall

Not fair! high precision may discount low recall or vice versa

# Average of Precision and Recall: A Single Metric

$$\text{Harmonic AVG-PR} = \frac{2}{\frac{1}{P} + \frac{1}{R}} = 2 \left( \frac{P \times R}{P + R} \right)$$

Same weights to Precision and Recall

Conservative! More toward the lower number.

$$\text{HarmonicMean}(a_1, a_2, a_3, a_4, \dots, a_n) = \frac{n}{\frac{1}{a_1} + \frac{1}{a_2} + \frac{1}{a_3} + \dots + \frac{1}{a_n}}$$



# $F_\beta$ -Measure

Weights to Precision and Recall Separately

$$F = \frac{1}{\alpha \frac{1}{P} + (1 - \alpha) \frac{1}{R}} \quad \text{or} \quad \left( \text{with } \beta^2 = \frac{1 - \alpha}{\alpha} \right) \quad F = \frac{(\beta^2 + 1)PR}{\beta^2 P + R}$$

$$F_\beta = \frac{(\beta^2 + 1)PR}{\beta^2 P + R} \quad \left\{ \begin{array}{ll} \beta > 1: & \text{favors Recall} \\ \beta = 1: & 2 \left( \frac{P \times R}{P + R} \right) \\ 0 \leq \beta < 1: & \text{favors Precision} \end{array} \right.$$

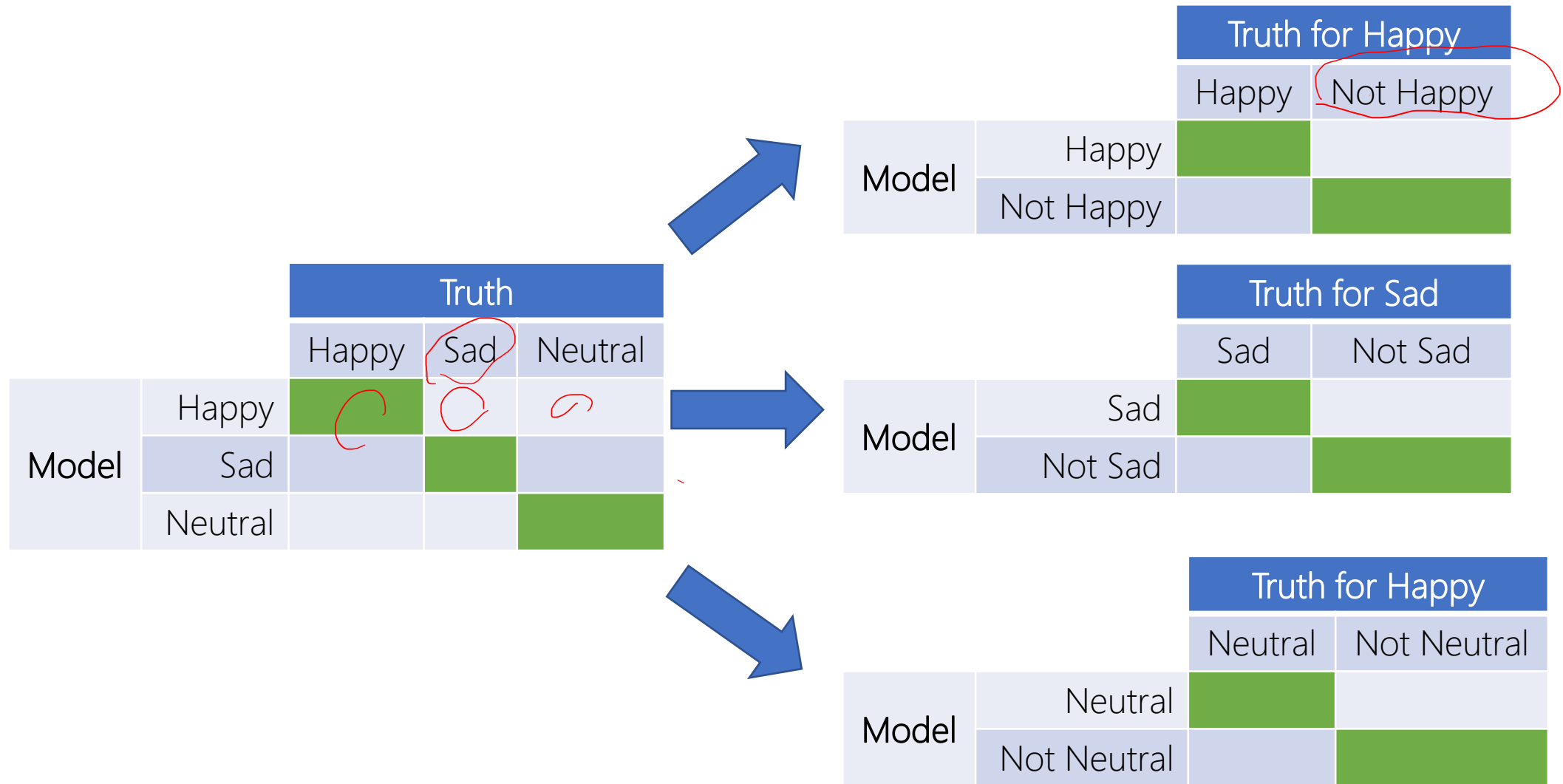


# Multiclass Evaluation

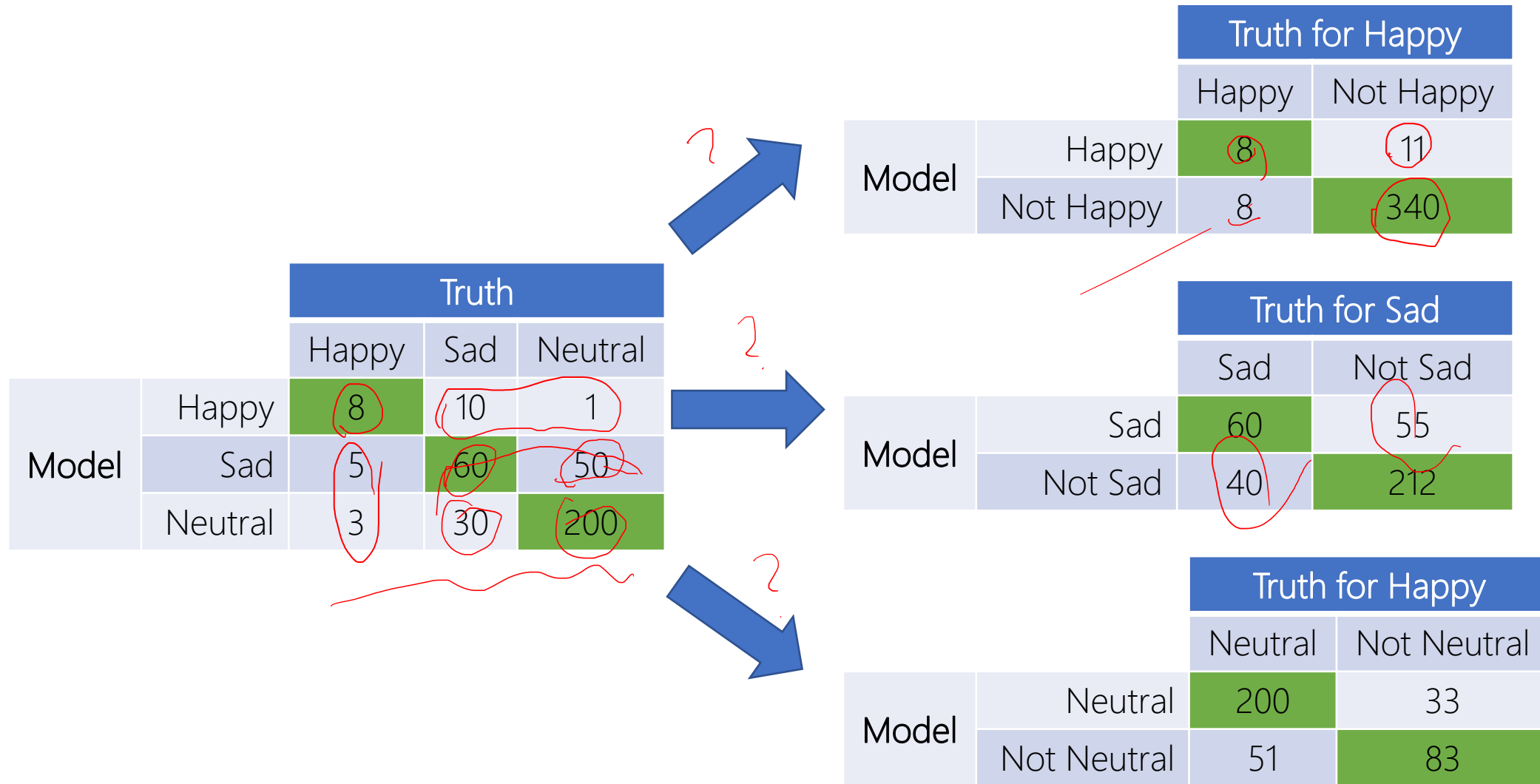
---



# Multiclass Evaluation




# Multiclass Evaluation






# Multiclass Evaluation

		Truth		
		Happy	Sad	Neutral
Model	Happy	8	10	1
	Sad	5	60	50
	Neutral	3	30	200




		Truth for Happy	
		Happy	Not Happy
Model	Happy	8	11
	Not Happy	8	340

$$P_{\text{happy}} = \frac{8}{8+11} = 0.42$$



		Truth for Sad	
		Sad	Not Sad
Model	Sad	60	55
	Not Sad	40	212

$$P_{\text{sad}} = \frac{60}{60+55} = 0.52$$



		Truth for Happy	
		Neutral	Not Neutral
Model	Neutral	200	33
	Not Neutral	51	83

$$P_{\text{neutral}} = \frac{200}{200+33} = 0.85$$

# Multiclass Evaluation

		Truth		
		Happy	Sad	Neutral
Model	Happy	8	10	1
	Sad	5	60	50
	Neutral	3	30	200

$$R_{\text{happy}} = \frac{8}{8+5+3} = ?$$

$$R_{\text{sad}} = \frac{60}{10+60+30} = ?$$

$$R_{\text{neutral}} = \frac{200}{1+50+200} = ?$$

$$P_{\text{happy}} = \frac{8}{8+10+1} = 0.42$$

$$P_{\text{sad}} = \frac{60}{5+60+50} = 0.52$$

$$P_{\text{neutral}} = \frac{200}{3+30+200} = 0.85$$



# Multiclass Evaluation: Macro-Avg

$$Macroavg = \frac{1}{K} \sum_{i=1}^K Metric_K$$

$$Macroavg = \frac{1}{3} [P_{happy} + P_{sad} + P_{neutral}]$$

The diagram illustrates the calculation of Macro-Avg precision for a 3-class problem. It starts with a confusion matrix and decomposes it into three 2-class problems, each with its own precision calculation.

**Confusion Matrix:**

		Truth		
		Happy	Sad	Neutral
Model	Happy	8	10	1
	Sad	5	60	50
	Neutral	3	30	200

**Truth for Happy (2-class):**

		Happy	Not Happy
Model	Happy	8	11
	Not Happy	8	340

$P_{happy} = \frac{8}{8+11} = 0.42$

**Truth for Sad (2-class):**

		Sad	Not Sad
Model	Sad	60	55
	Not Sad	40	212

$P_{sad} = \frac{60}{60+55} = 0.52$

**Truth for Neutral (2-class):**

		Neutral	Not Neutral
Model	Neutral	200	33
	Not Neutral	51	83

$P_{neutral} = \frac{200}{200+33} = 0.85$

A red bracket on the right groups the three precision values ( $P_{happy}$ ,  $P_{sad}$ ,  $P_{neutral}$ ) which are then averaged to find the Macro-Avg.

# Multiclass Evaluation: Micro-Avg

		Truth		
		Happy	Sad	Neutral
Model	Happy	8	10	1
	Sad	5	60	50
	Neutral	3	30	200

Pool



		True	Not True
		268	99
Model	True	268	99
	Not True	99	?

$R_{\text{true}} = ?$

$P_{\text{true}} = ?$





# Macro vs. Micro Averaging

---

# Macro vs. Micro

---

Micro is dominated by the more frequent class.

Micro better reflects the statistics of the smaller classes.

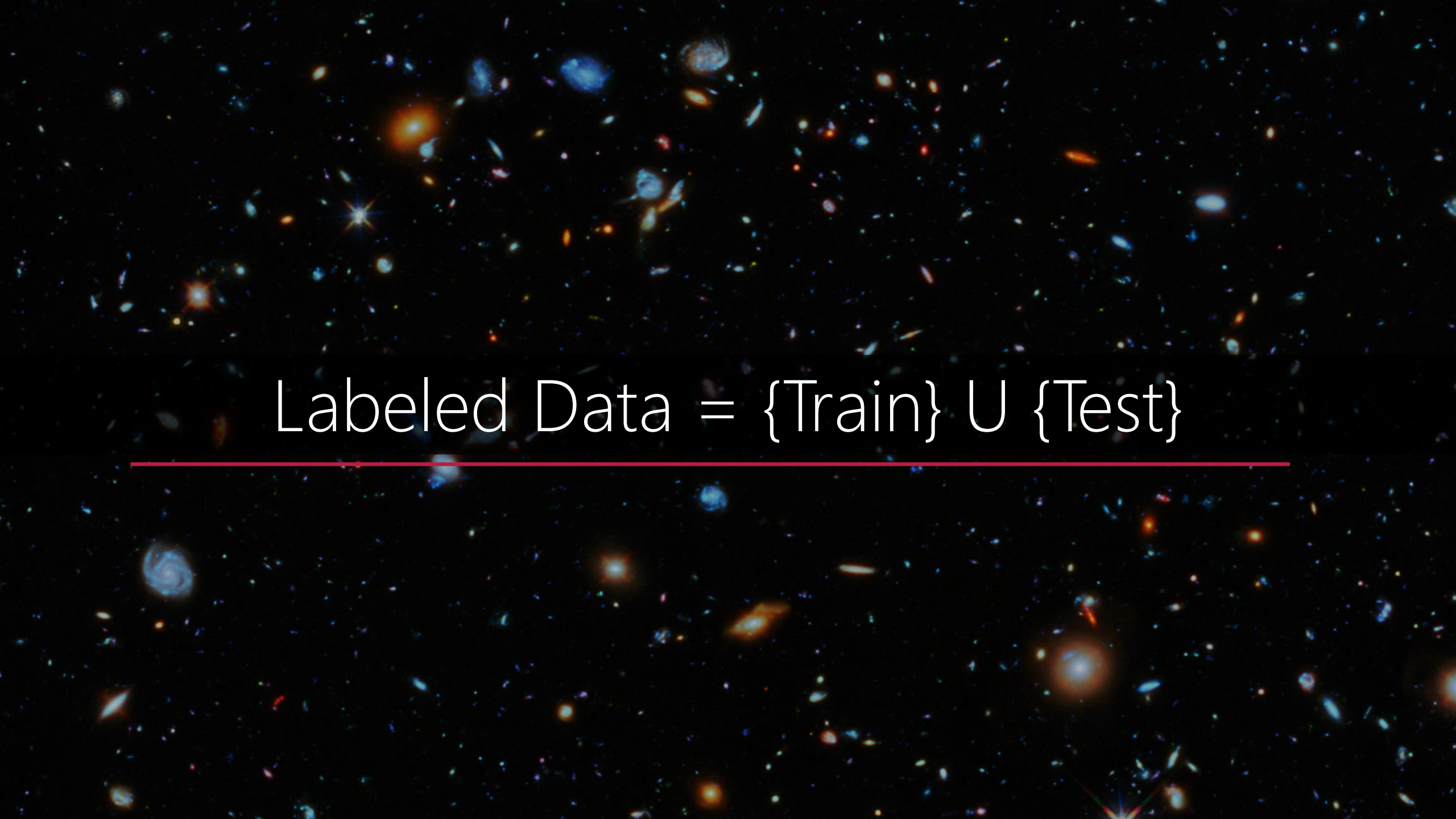
Micro more appropriate when performance on all the classes is equally important.





Calculate the Metrics

---



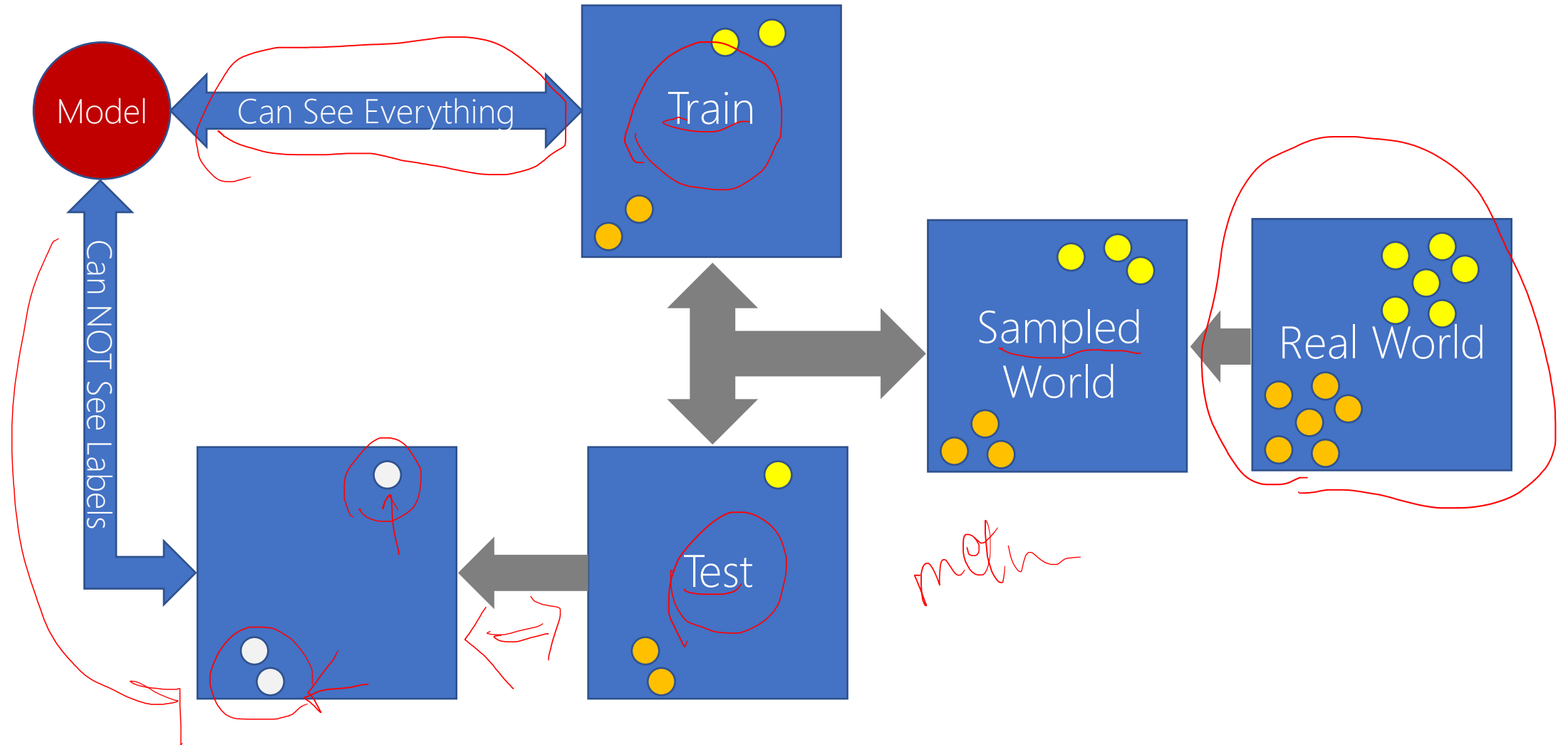
Labeled Data = {Train} U {Test}

---



# Labeled Data = {Train} U {Test}

---



Labeled Data = {Train} U {Test}

---

$$\{Train\} \cap \{Test\} \stackrel{?}{=} \emptyset$$



Imbalance Labeled Data = {Train} U {Test}

---

$$\{Train\} \cap \{Test\} \stackrel{?}{=} \emptyset$$

Train and test sets presumably follow same distribution!





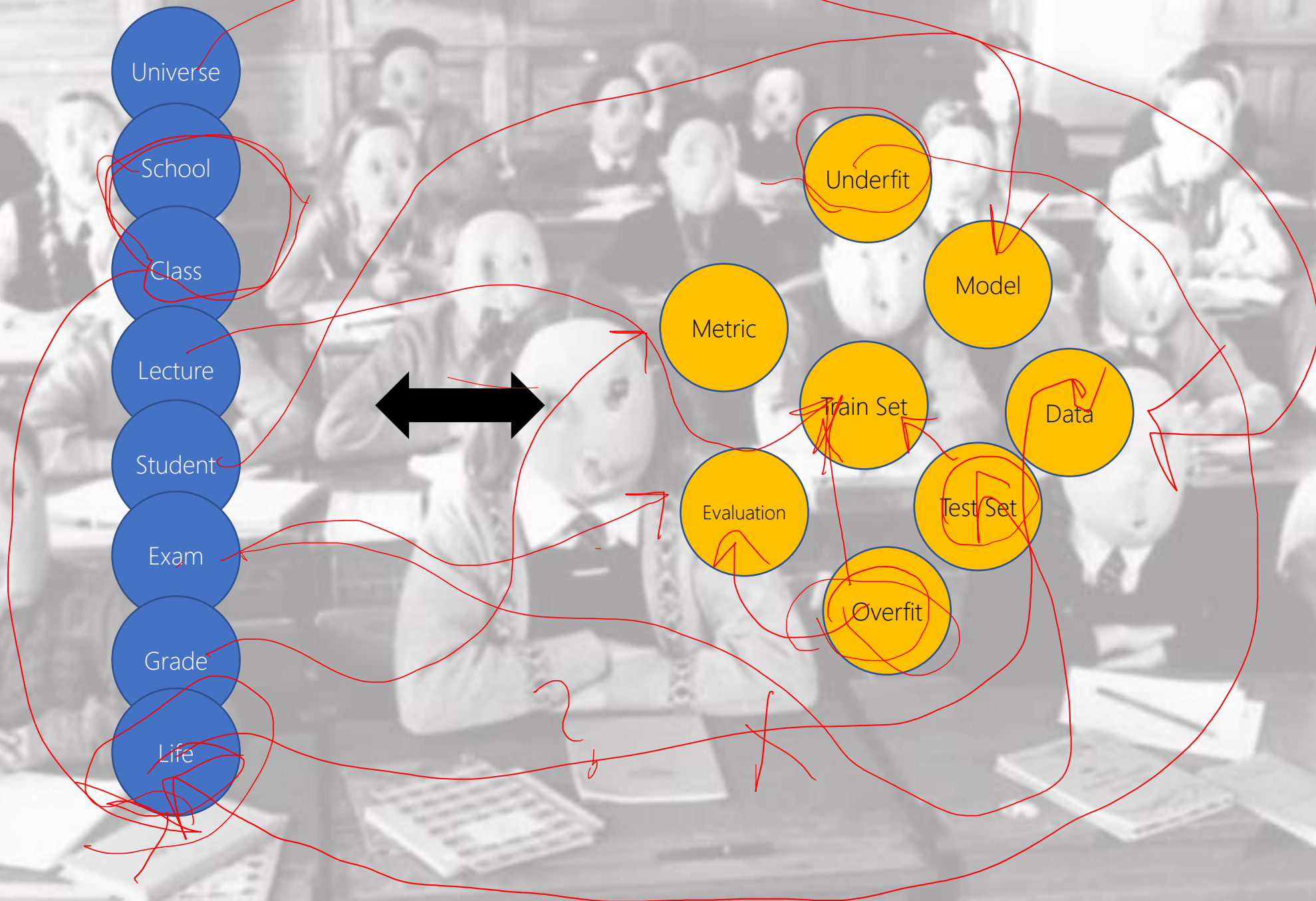
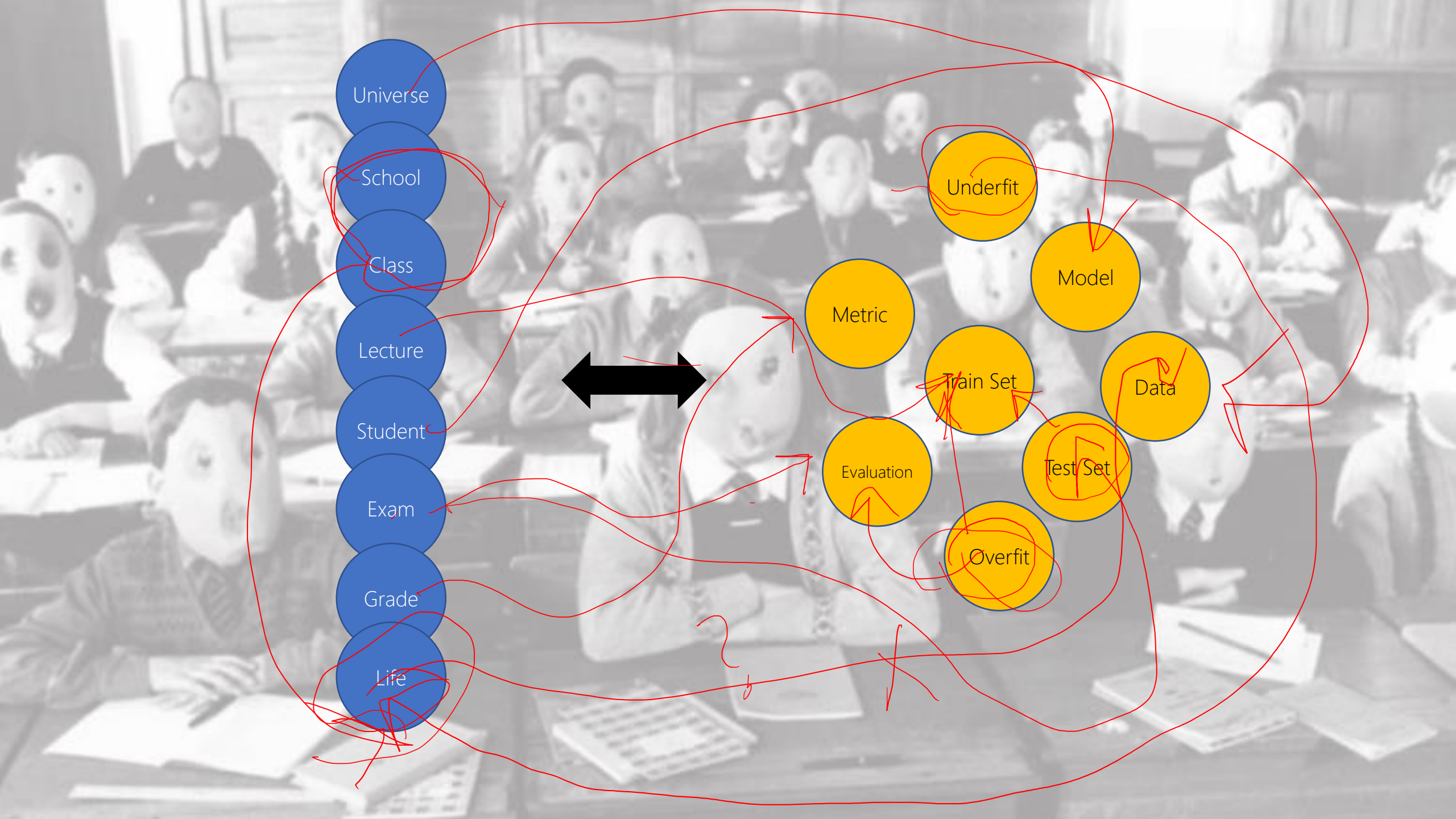
Underfit  $\rightarrow$  Balance fit  $\leftarrow$  Overfit

---



Another Brick in the Wall  
- Pink Floyd  
- The Wall







---

# Model Tuning

---

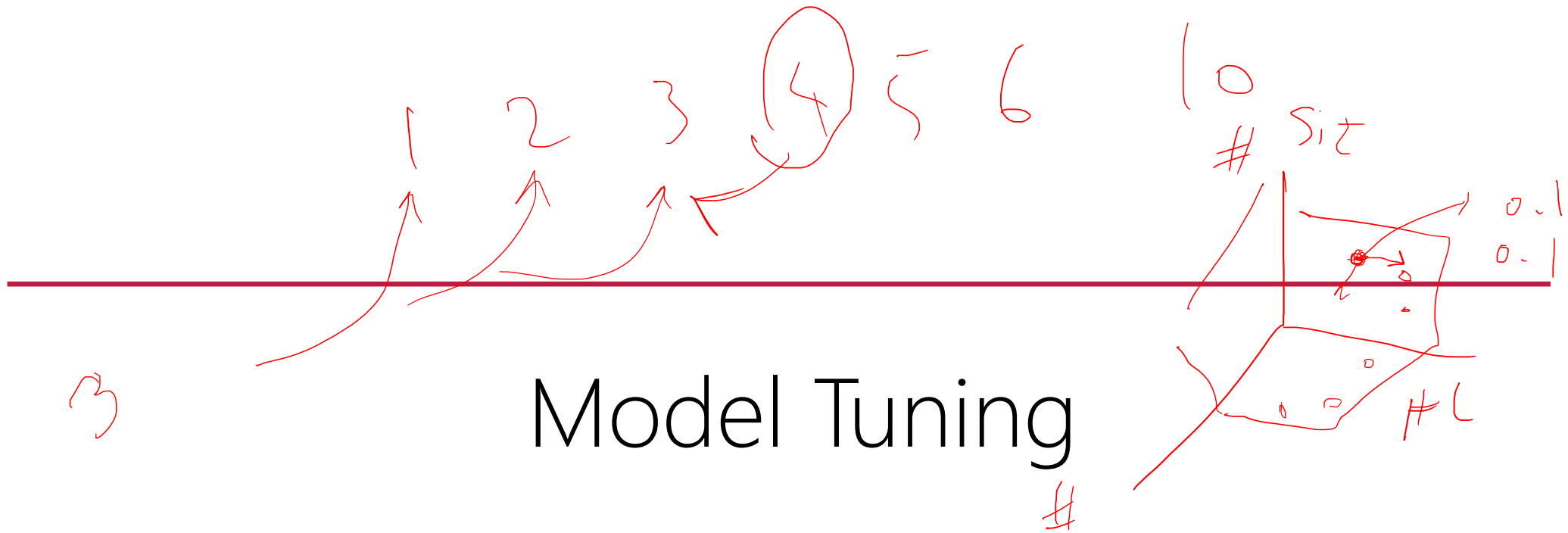
Find the best running settings of the mode

- n-Gram LM: what is the best n? hyper
- Probs. assumptions
- Neural Model: #layers, Activation functions



## Checking the performance of model on Train and Test

Blind grid search! Brute-force



## Model Tuning

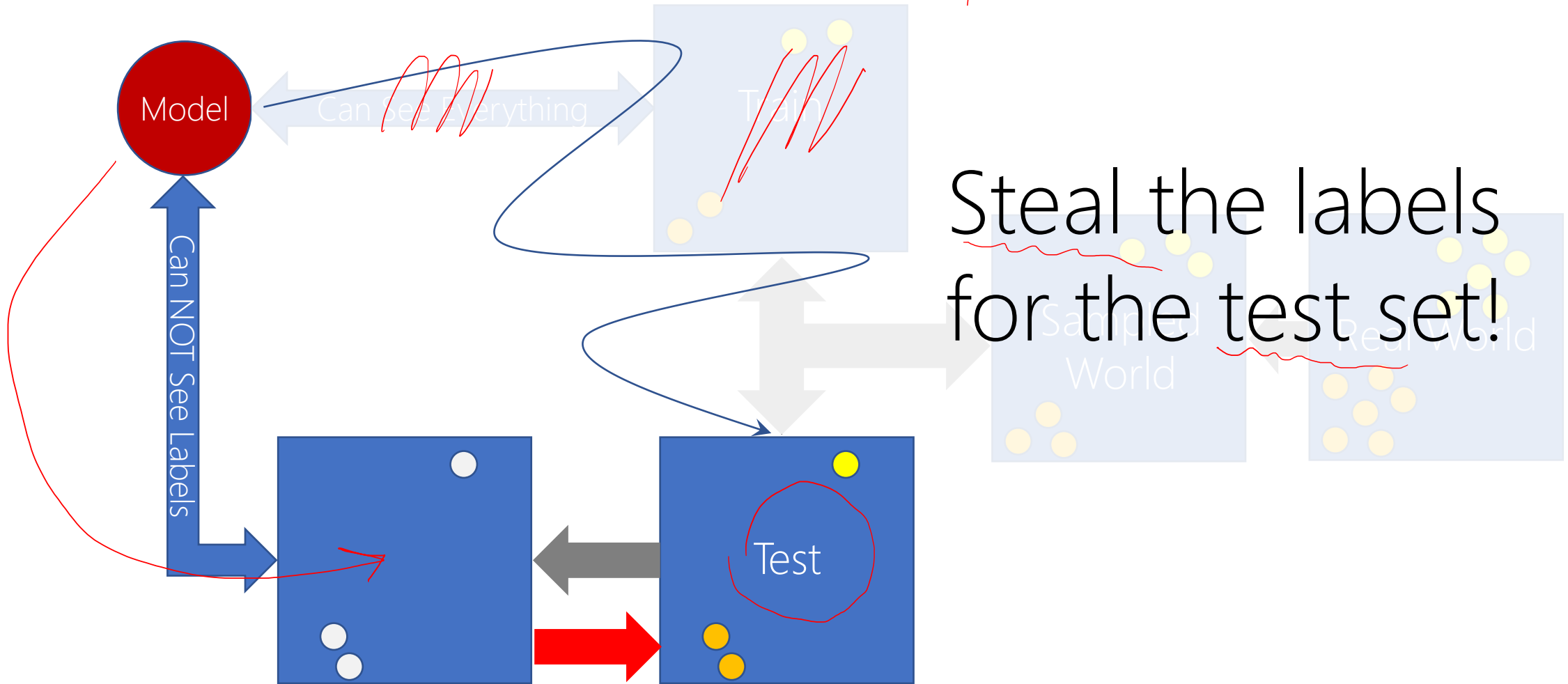
Find the best running settings of the mode

- Learn the performance of model on Train and Test
- For all different possibilities

Guided grid search!




$$\text{Labeled Data} = \{\text{Train}\} \cup \{\text{Test}\}$$



🏆

Featured Code Competition

 PetFinder.my · 2,023 teams · 2 years ago

# PetFinder.my Adoption Prediction

How cute is that doggy in the shelter?

\$25,000

Prize Money

Overview

Data


Code

Discussion

Leaderboard

Rules


New Topic



Mongrel Jedi

## PetFinder.my Contest: 1st Place Winner Disqualified

Posted in [petfinder-adoption-prediction](#) a year ago



307

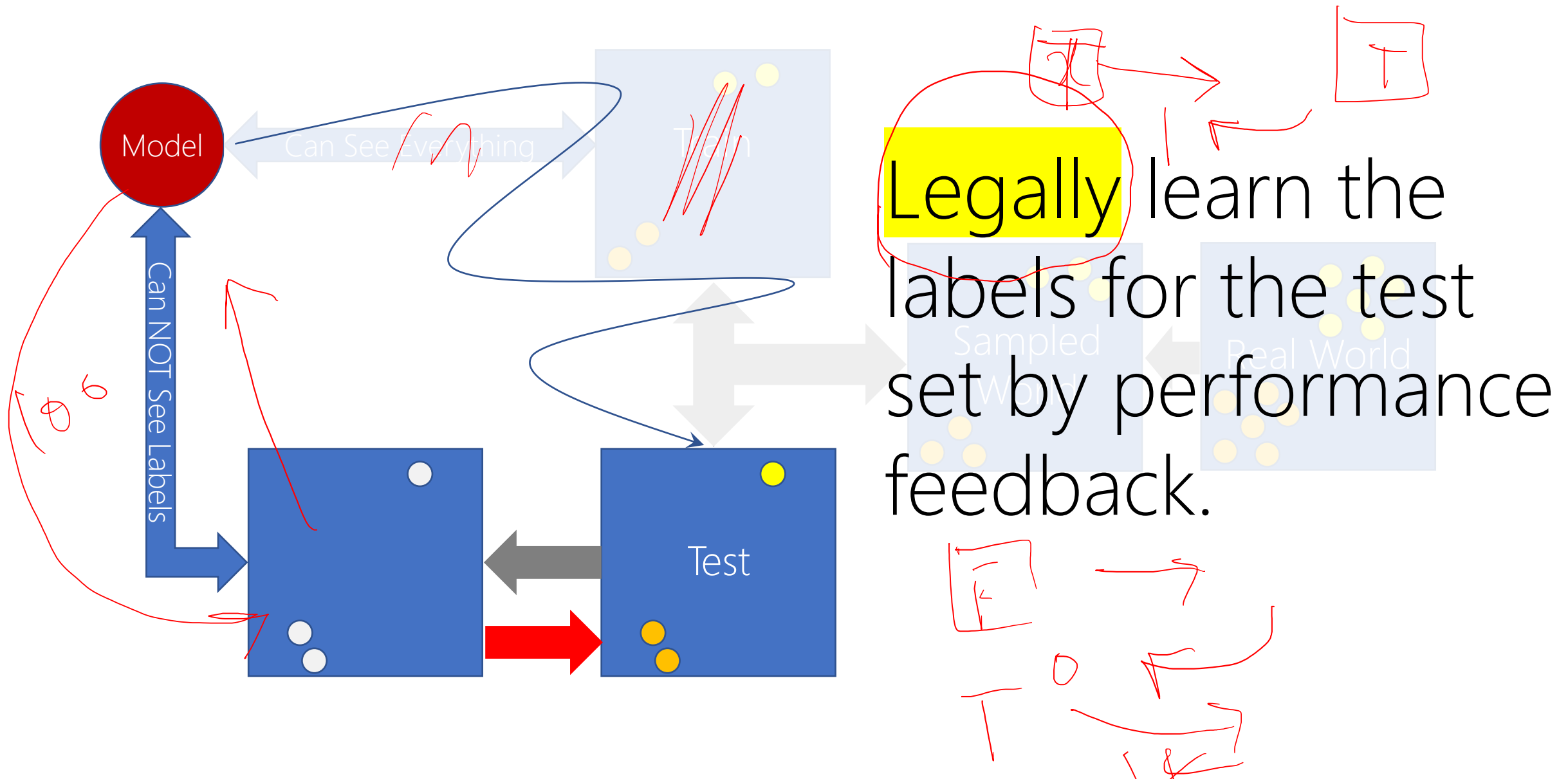
Dear Participants,

We would like to announce that the 1st Place Team, Bestpetting has been disqualified from the contest for cheating. The Kaggle Grandmaster cheater has also been permanently banned on this platform as the evidence points towards him being the key party behind this fraudulent activity.

Here is what the Bestpetting team did in the [PetFinder.my](#) contest:

- They fraudulently obtained adoption speed answers for the private test data (possibly by scraping our website)

$$\text{Labeled Data} = \{\text{Train}\} \cup \{\text{Test}\}$$





$$\text{Labeled Data} = \{\text{Train}\} \cup \{\text{Test}\}$$

---



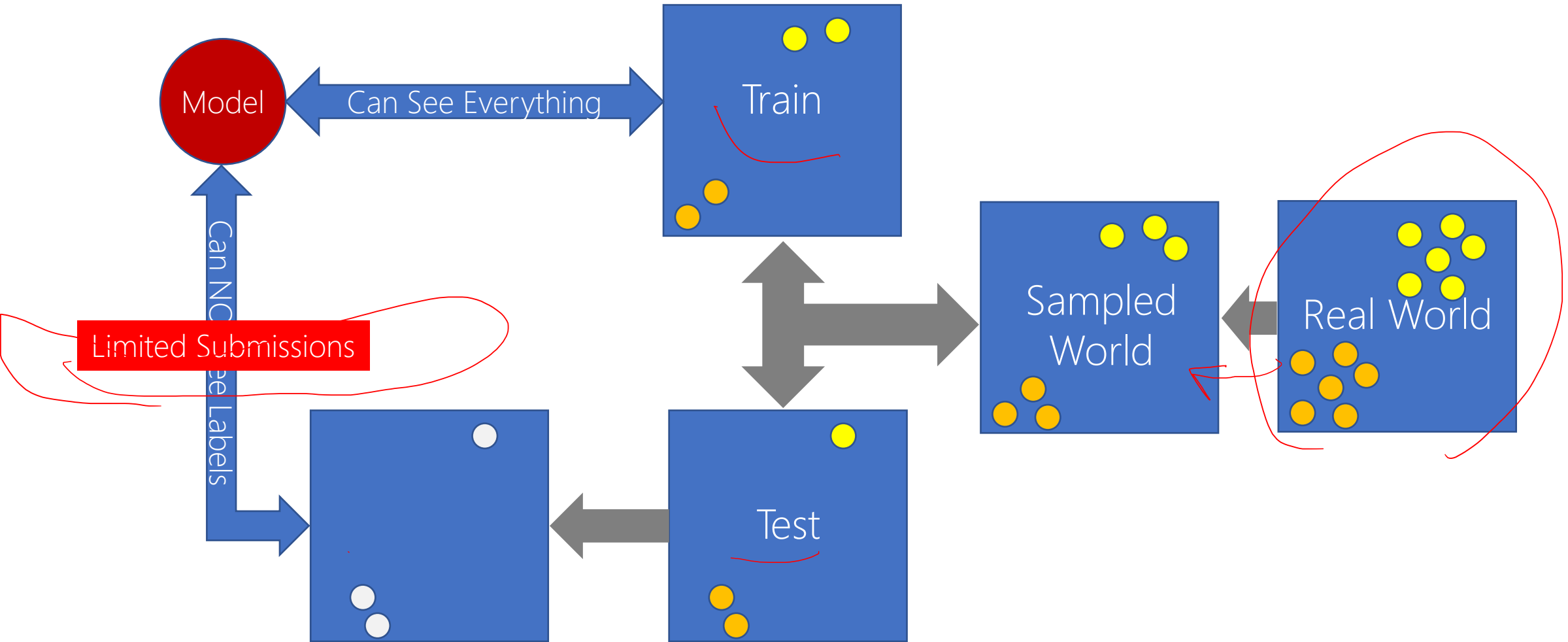


Labeled Data =  $\{\{\text{Train}\} \cup \{\text{Valid}\}\} \cup \{\text{Test}\}$

The model intentionally ignores parts of his available knowledge and challenges itself to uncover those parts!

# Labeled Data = {Train} U {Test}

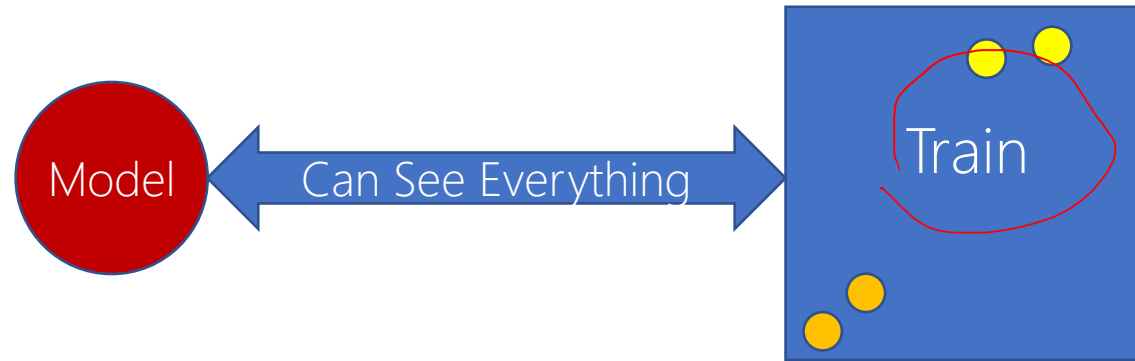
---



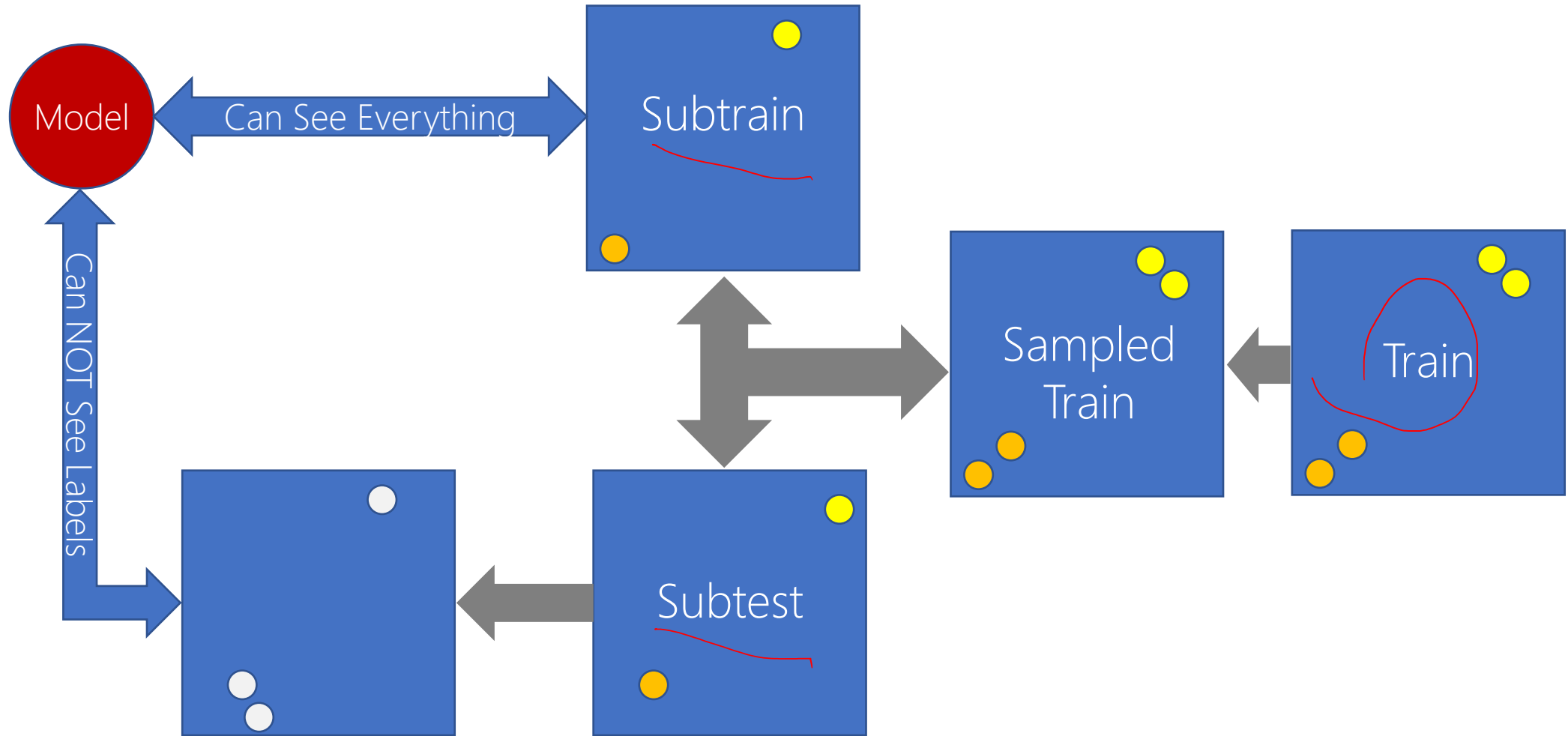


# Labeled Data = {Train}

---

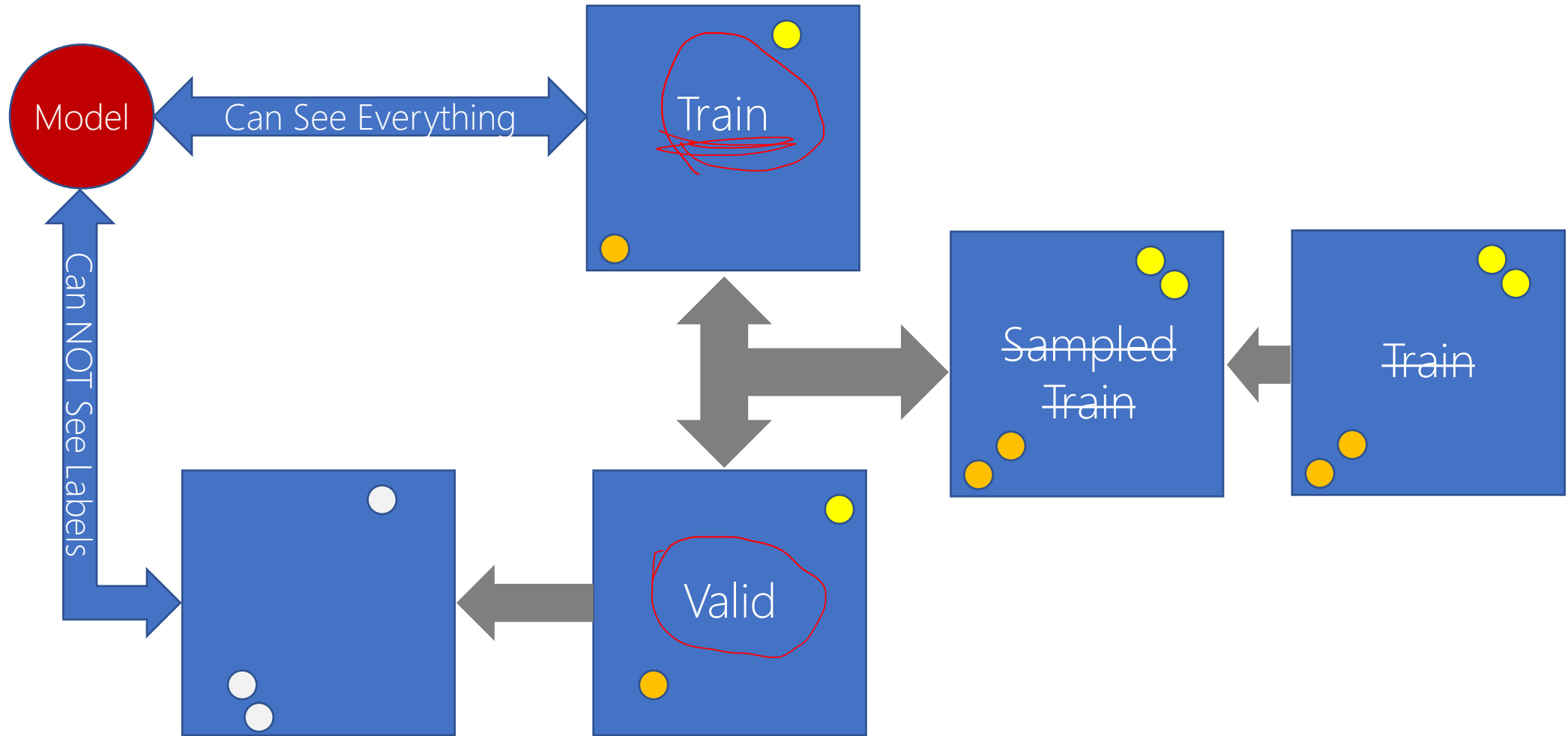


Labeled Data = {{Subtrain} U {Subtest}}



# Labeled Data = $\{\text{Train}\} \cup \{\text{Valid}\}$

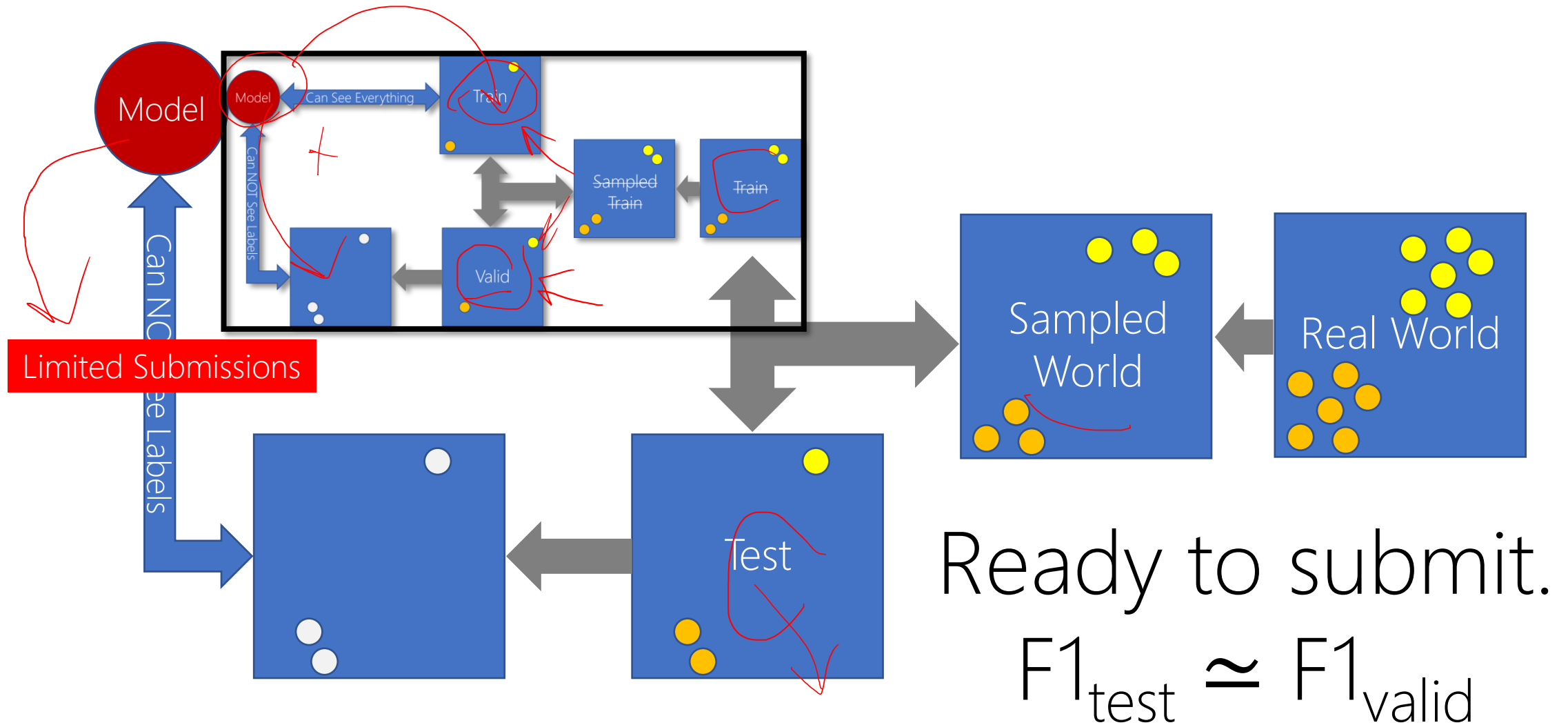
---





# Labeled Data = {Train} U {Test}

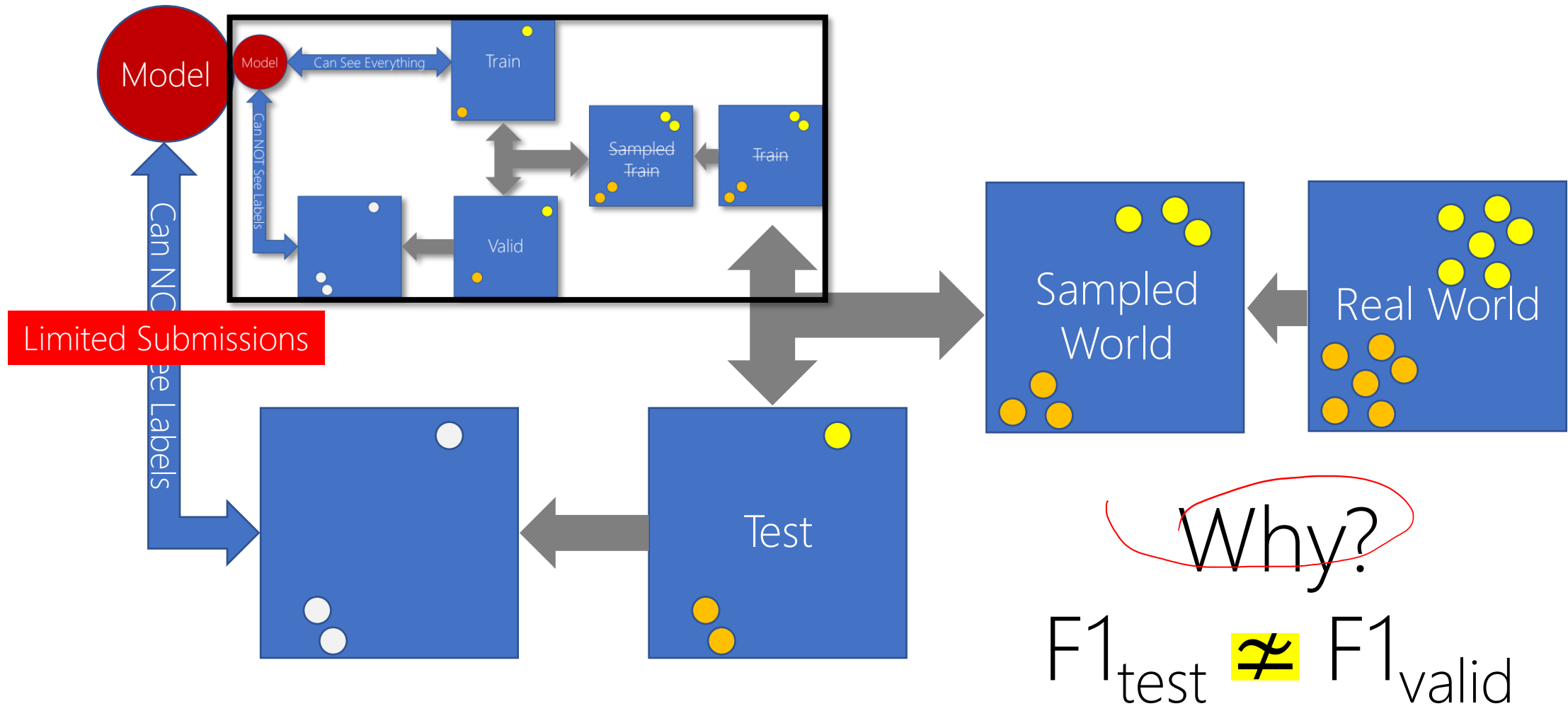
---





# Labeled Data = {Train} U {Test}

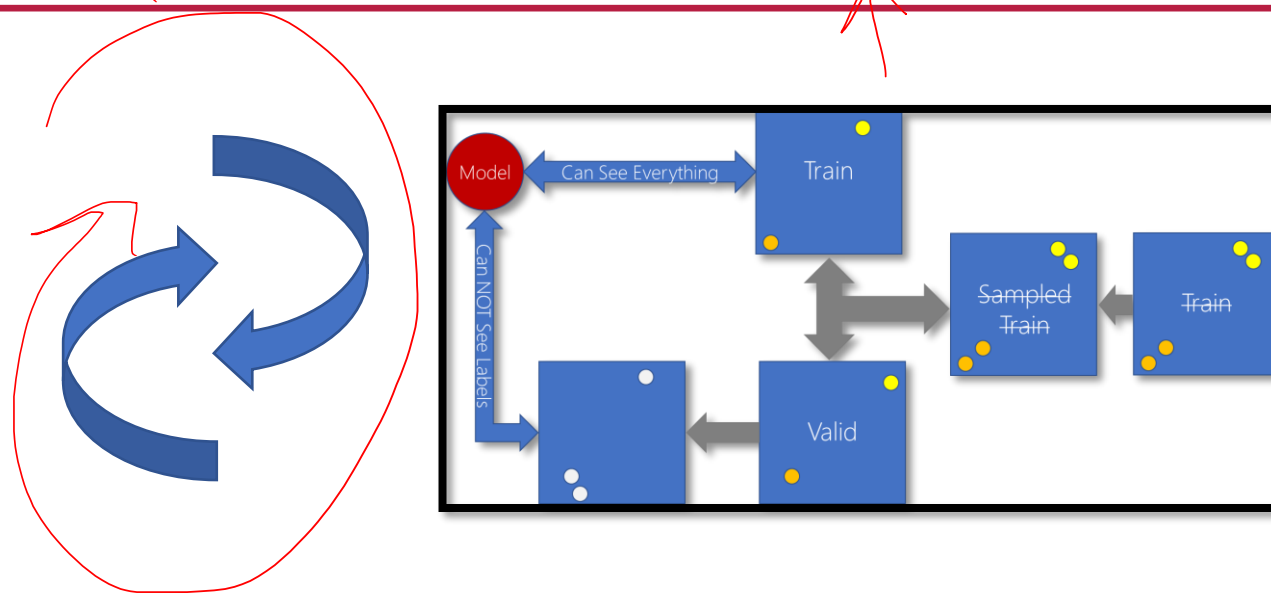
---



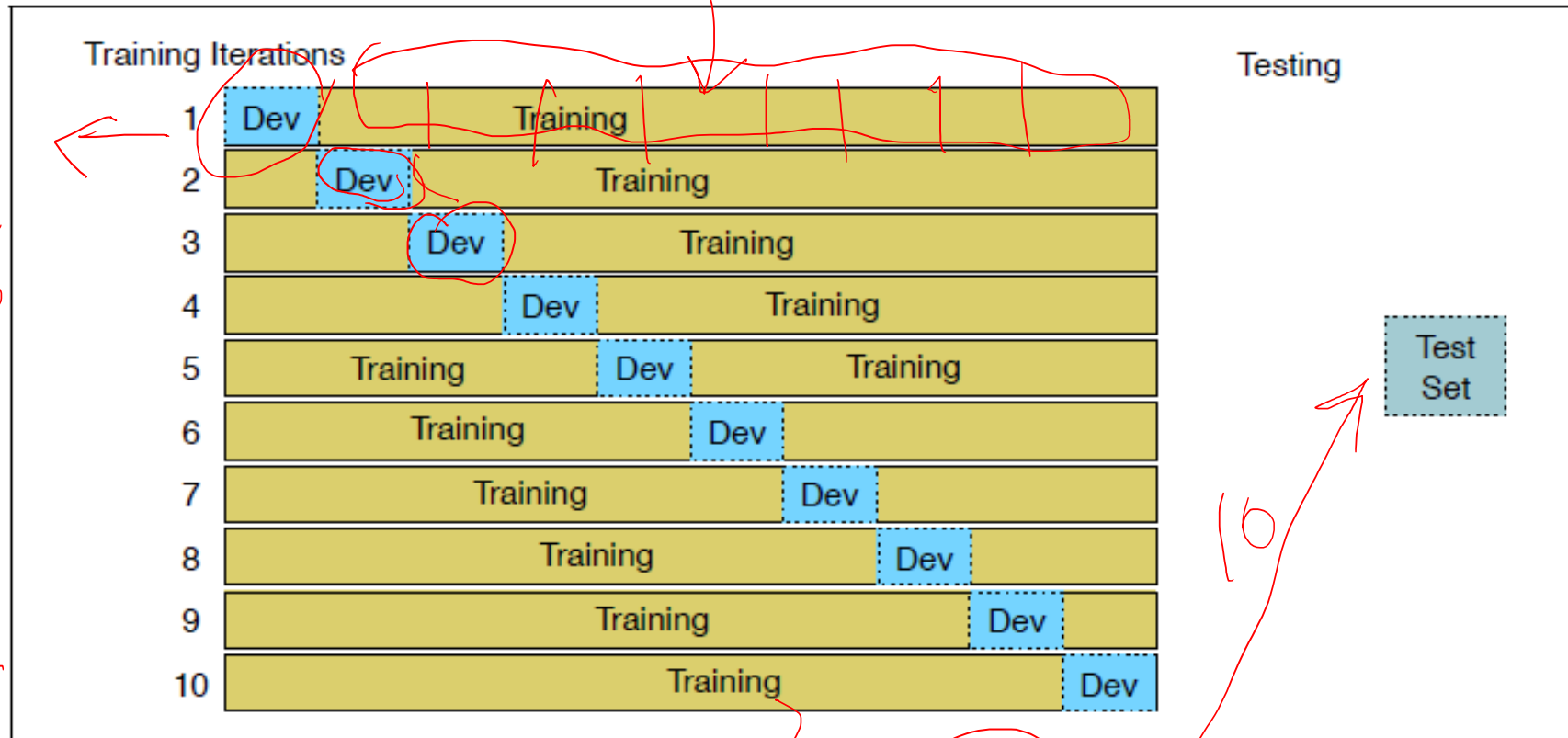


# Cross-Validation

1 practice vs. Multiple practice

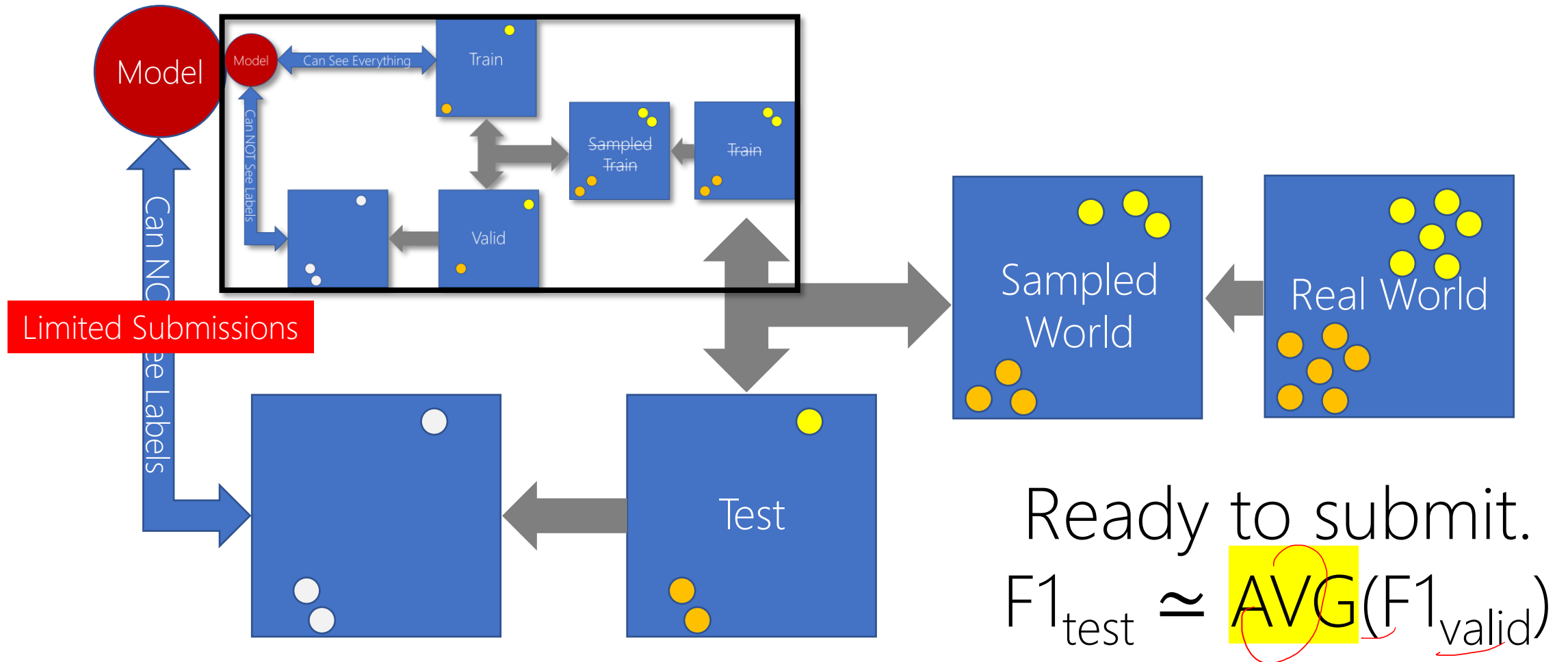


# Cross-Validation



**Figure 4.7** 10-fold cross-validation

$$\text{Labeled Data} = \{\text{Train}\} \cup \{\text{Test}\}$$





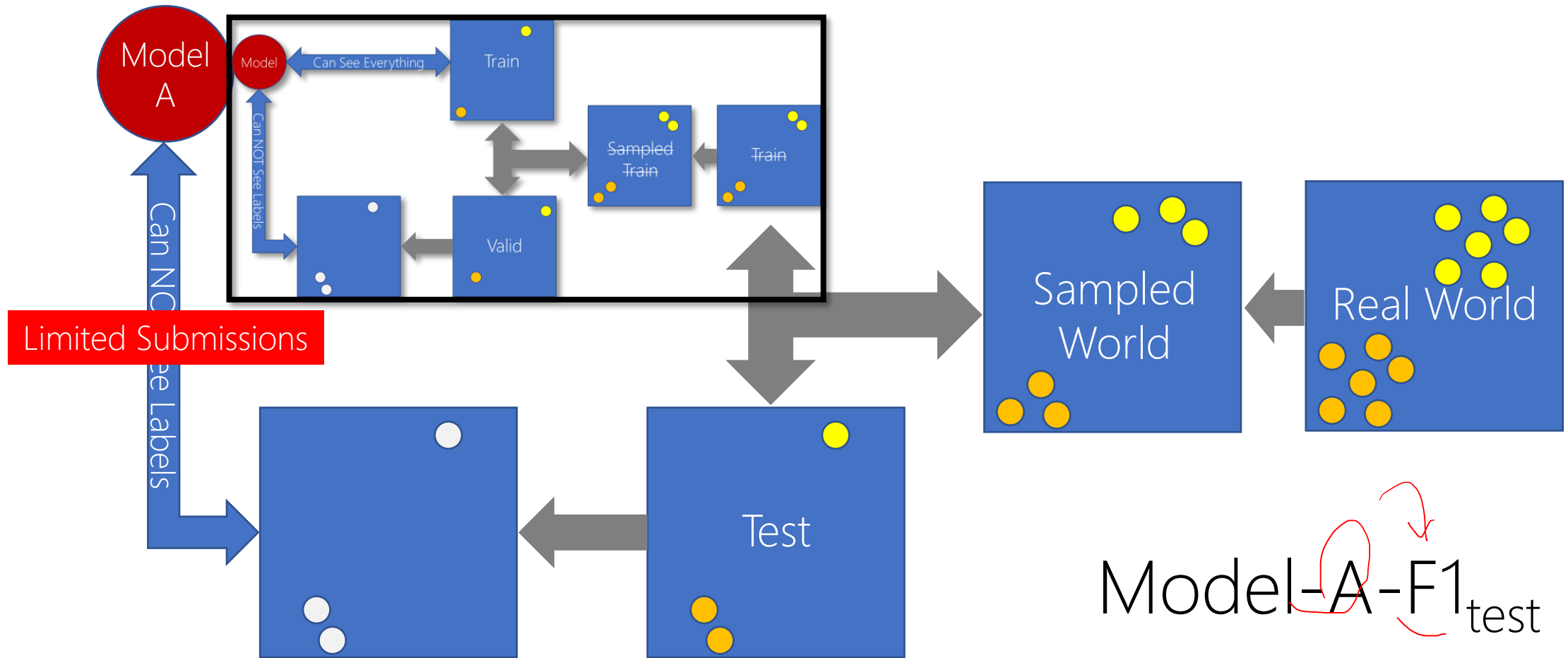
---

# Statistical Significance Testing

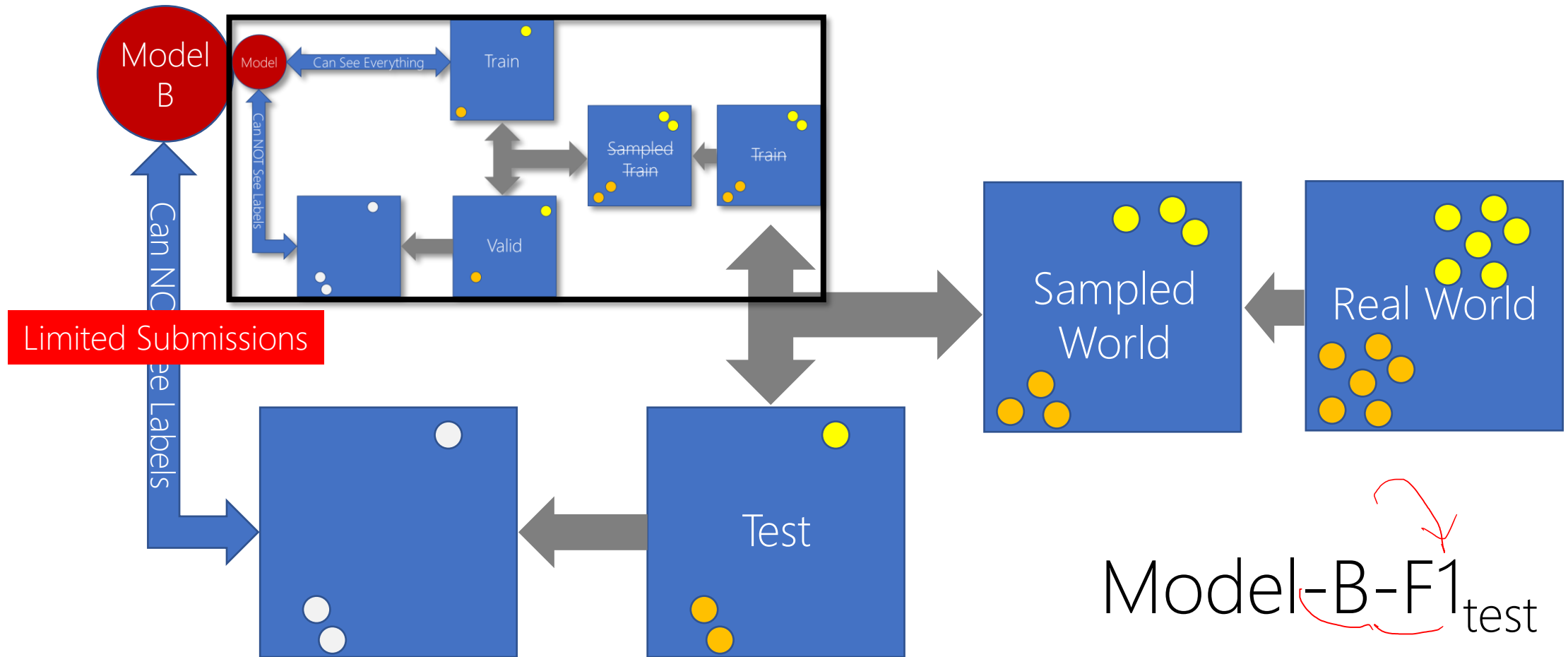
---

## Model Comparison

# Statistical Significance Testing



# Statistical Significance Testing





# Statistical Significance Testing

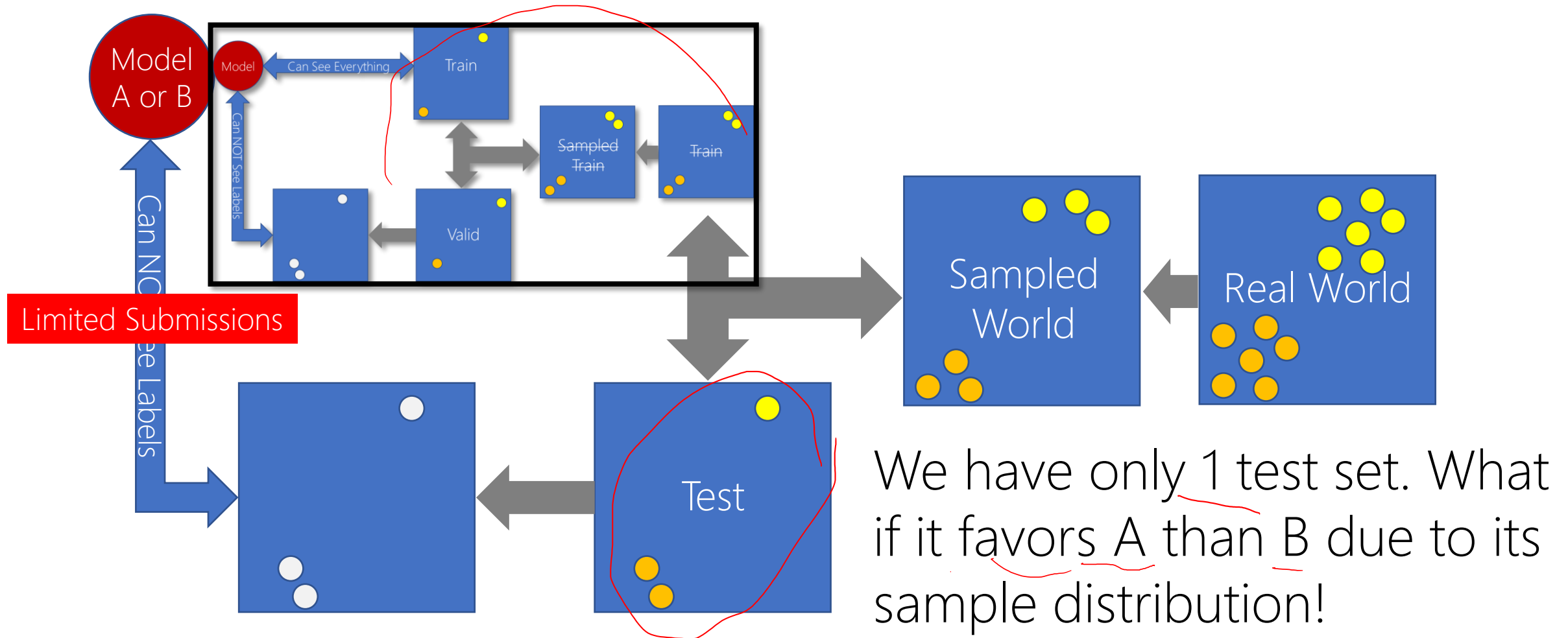
---

Test Sets	A-F1	E.g.,	B-F1	E.g.,
→ Test-1	A-F1 <sub>test-1</sub>	0.99	B-F1 <sub>test-1</sub>	0.6

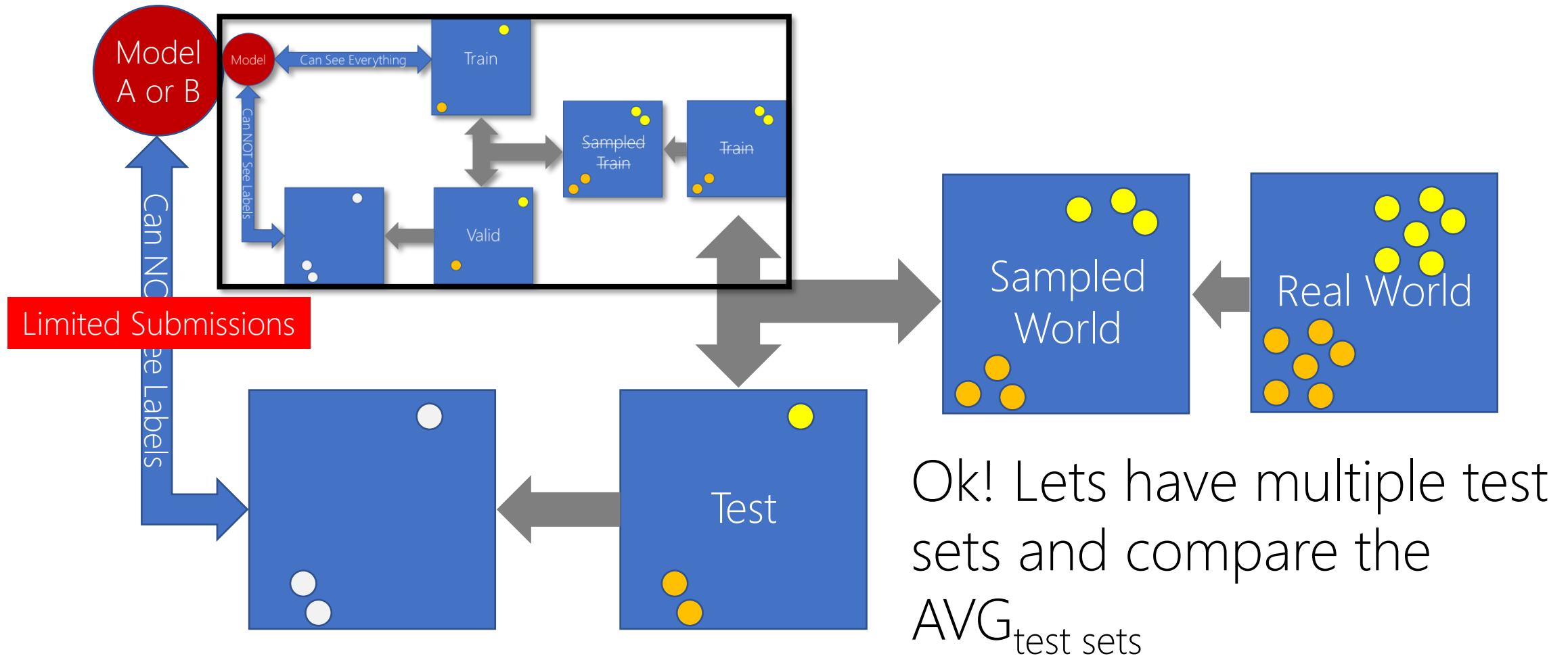
A is better than B.

A is significantly better than B →  $0.99 \gg 0.6$

# Statistical Significance Testing



# Statistical Significance Testing





# Statistical Significance Testing

Test Sets	A-F1	E.g.,	B-F1	E.g.,
Test-1	$A-F1_{\text{test-1}}$	0.99	$B-F1_{\text{test-1}}$	0.6
Test-2	$A-F1_{\text{test-2}}$	0.5	$B-F1_{\text{test-2}}$	0.6
Test-3	$A-F1_{\text{test-3}}$	0.5	$B-F1_{\text{test-3}}$	0.6
Test-4	$A-F1_{\text{test-4}}$	0.5	$B-F1_{\text{test-4}}$	0.6
AVG		0.6225		0.6

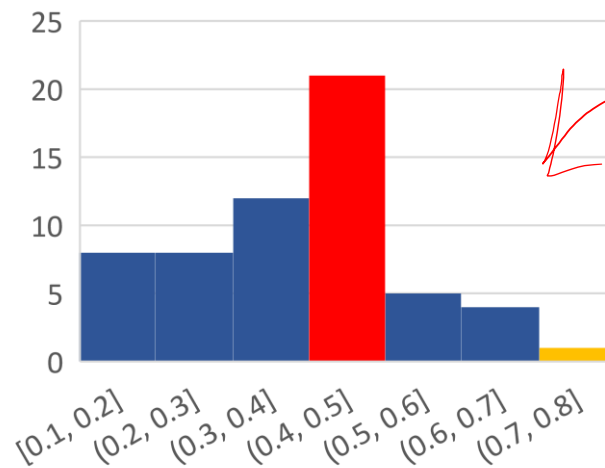
- 1) By average, A is better than B. However, clearly B is better than A.
- 2) By average, A is better than B but only slightly NOT significantly!

What is the problem here?

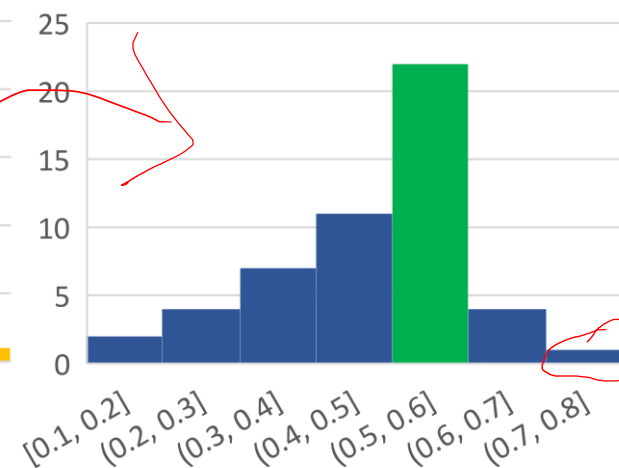
# Statistical Significance Testing

Test Sets	A-F1	E.g.,	B-F1	E.g.,
Test-1	$A-F1_{\text{test-1}}$	0.99	$B-F1_{\text{test-1}}$	0.6
Test-2	$A-F1_{\text{test-2}}$	0.5	$B-F1_{\text{test-2}}$	0.6
Test-3	$A-F1_{\text{test-3}}$	0.5	$B-F1_{\text{test-3}}$	0.6
Test-4	$A-F1_{\text{test-4}}$	0.5	$B-F1_{\text{test-4}}$	0.6
AVG		0.6225		0.6

A-F1 Histogram

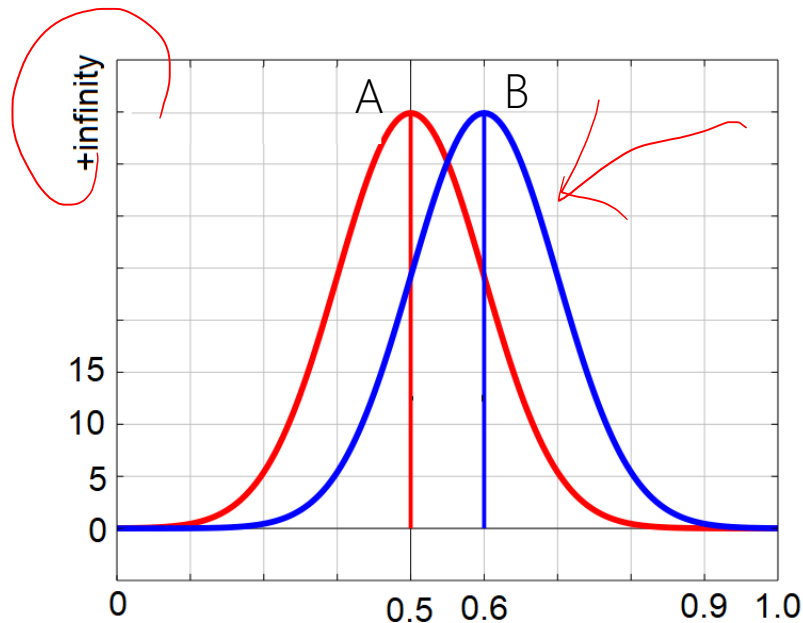


B-F1 Histogram



# Statistical Significance Testing

Test Sets	A-F1	E.g.,	B-F1	E.g.,
Test-1	$A-F1_{\text{test-1}}$	0.99	$B-F1_{\text{test-1}}$	0.6
Test-2	$A-F1_{\text{test-2}}$	0.5	$B-F1_{\text{test-2}}$	0.6
Test-3	$A-F1_{\text{test-3}}$	0.5	$B-F1_{\text{test-3}}$	0.6
Test-4	$A-F1_{\text{test-4}}$	0.5	$B-F1_{\text{test-4}}$	0.6
AVG		0.6225		0.6



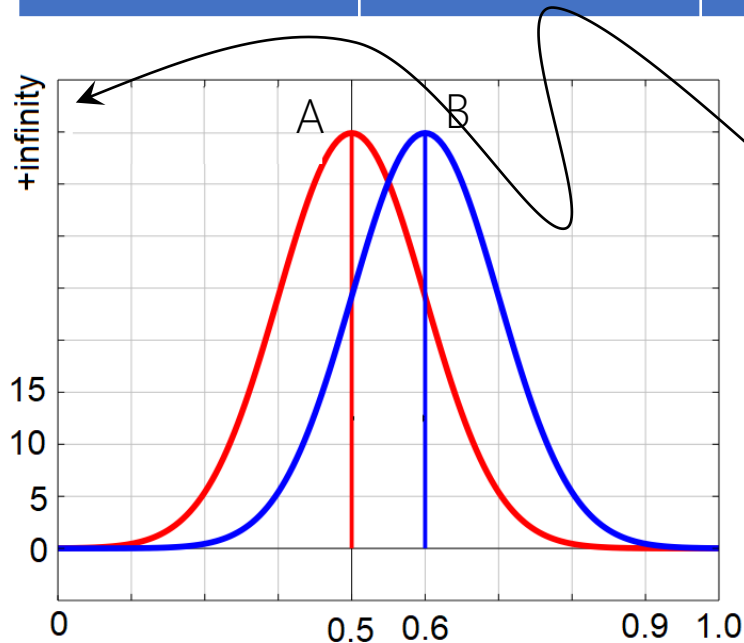
$$P(a < X < b) = \int_a^b f(x) dx.$$

$$P(0.2 < X < 0.4) > P(0.7 < X)$$

$$E(X) < E(Y)$$

# Statistical Significance Testing

Test Sets	A-F1	E.g.,	B-F1	E.g.,
Test-1	$A-F1_{\text{test-1}}$	0.99	$B-F1_{\text{test-1}}$	0.6
Test-2	$A-F1_{\text{test-2}}$	0.5	$B-F1_{\text{test-2}}$	0.6
Test-3	$A-F1_{\text{test-3}}$	0.5	$B-F1_{\text{test-3}}$	0.6
Test-4	$A-F1_{\text{test-4}}$	0.5	$B-F1_{\text{test-4}}$	0.6
AVG		0.6225		0.6



- 1) Labeled data is already expensive.
- 2) Sometimes testing is slow.

Reporting for a lot of runs on different test sets is very challenging!



A deep-field astronomical image showing a vast field of galaxies in various colors (blue, orange, red) against a black background. The galaxies are of different shapes and sizes, some appearing as bright, distinct objects while others are fainter and more distant.

# Statistical Significance Test

---

*t*-test



"I hope to arrive to my death,  
late,  
in love,  
and a little drunk."

— Atticus, anonymous Canadian poet