

# Machine Learning for Healthcare

Lessons From Across The Healthcare ML Lifecycle

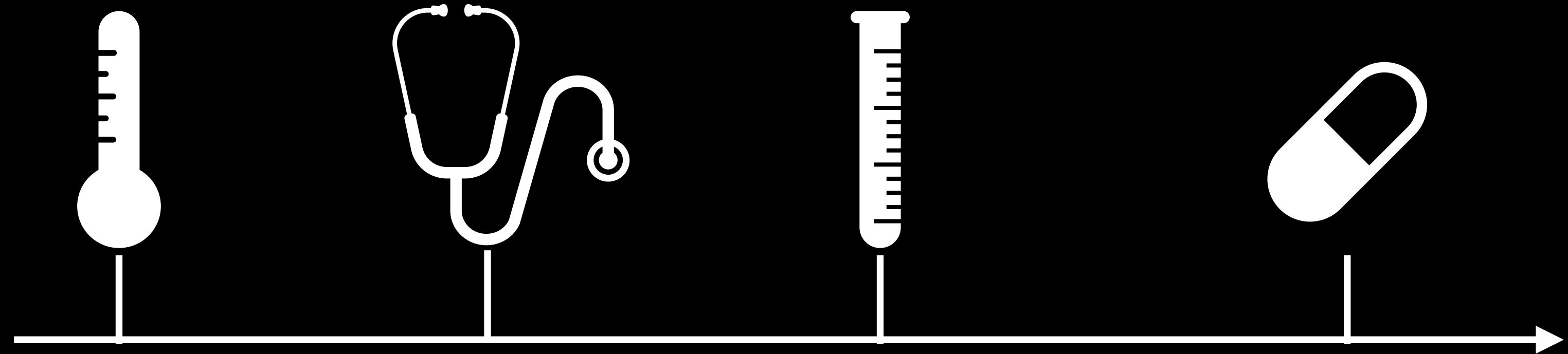
WPI Business Week  
Erkin Ötleş  
October 2023



# Time plays an important role in medicine.

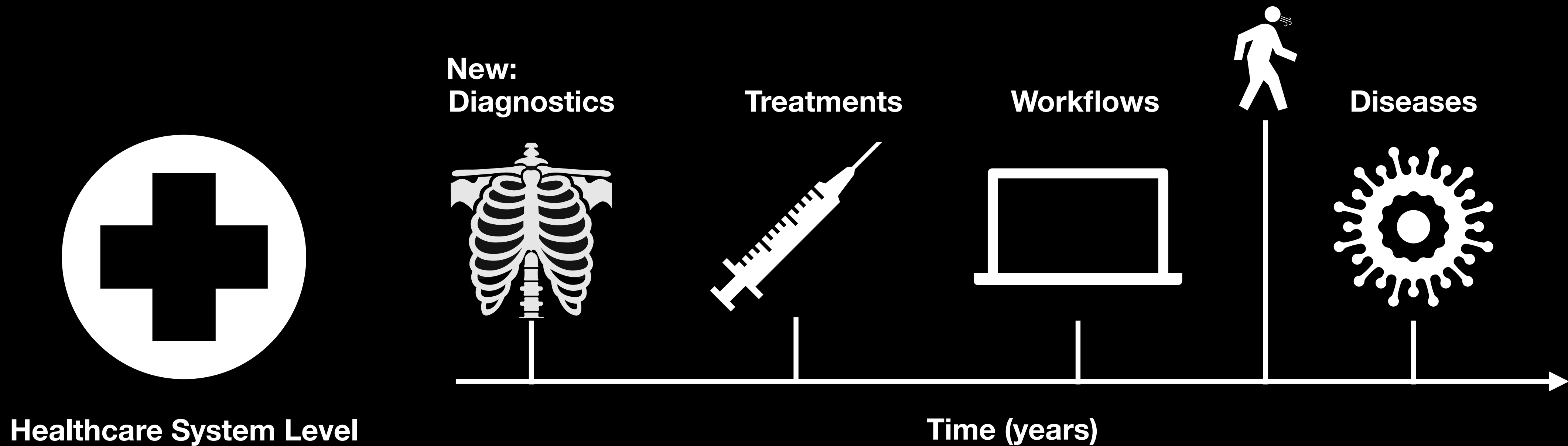


Patient Level



Time (days)

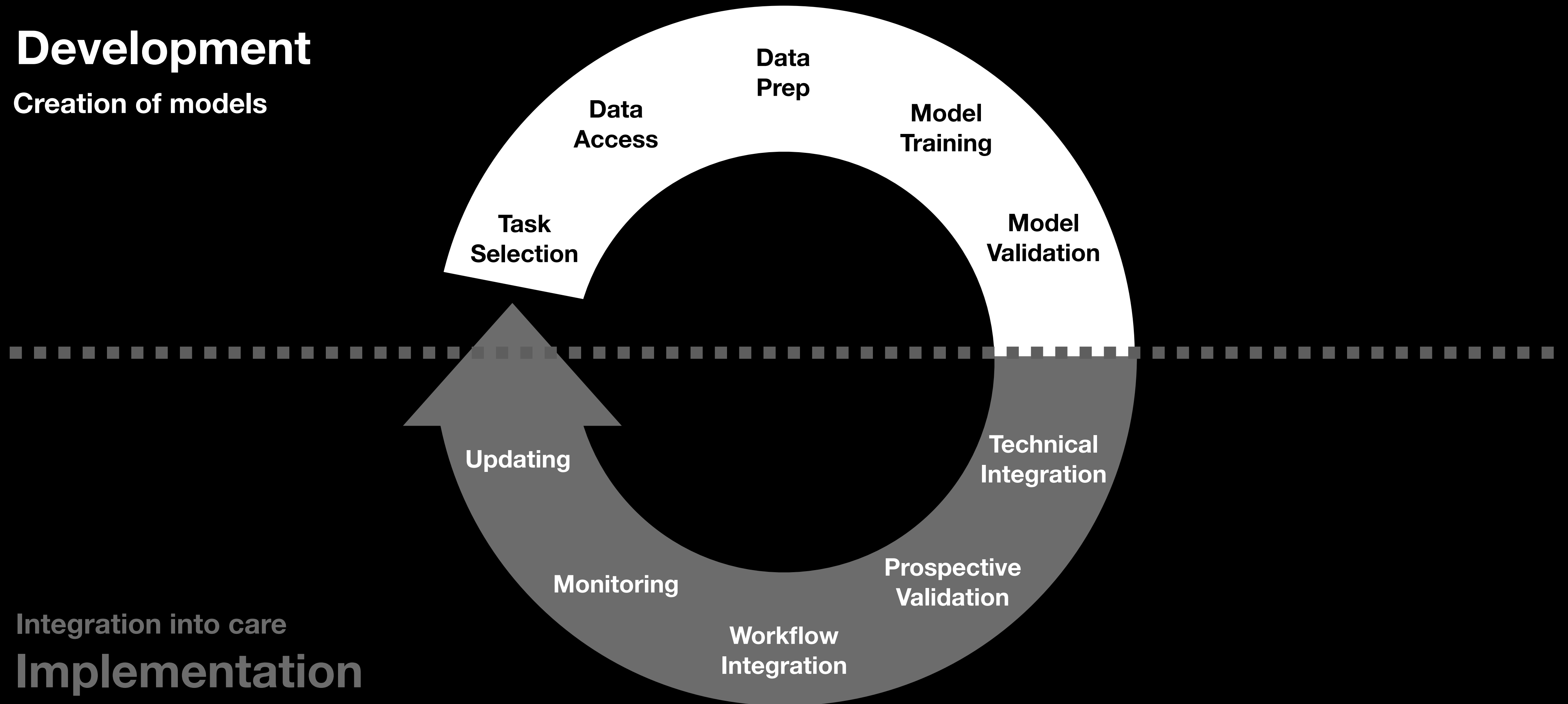
# Time plays an important role in medicine.



# Time complicates ML model development & implementation.

## Development

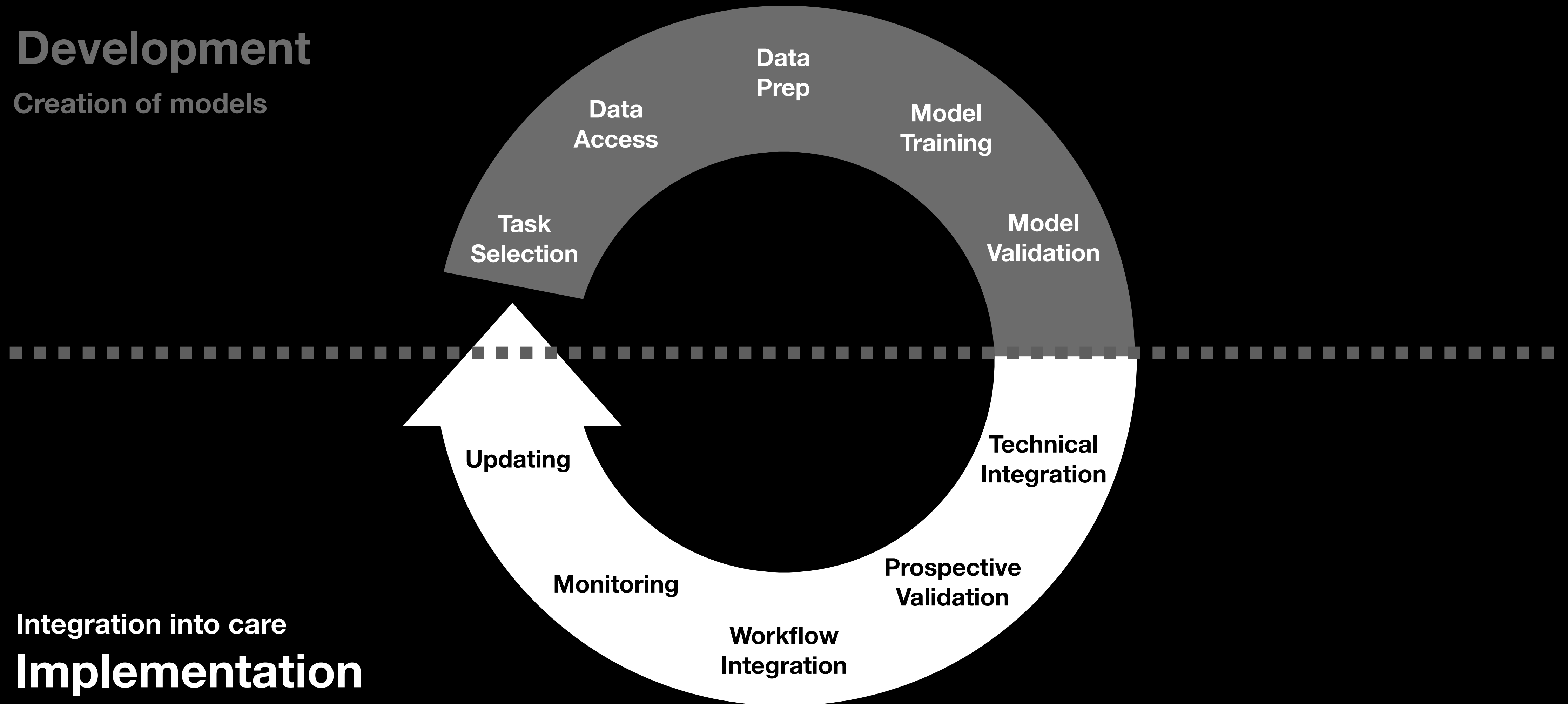
Creation of models



# Time complicates ML model development & implementation.

## Development

Creation of models



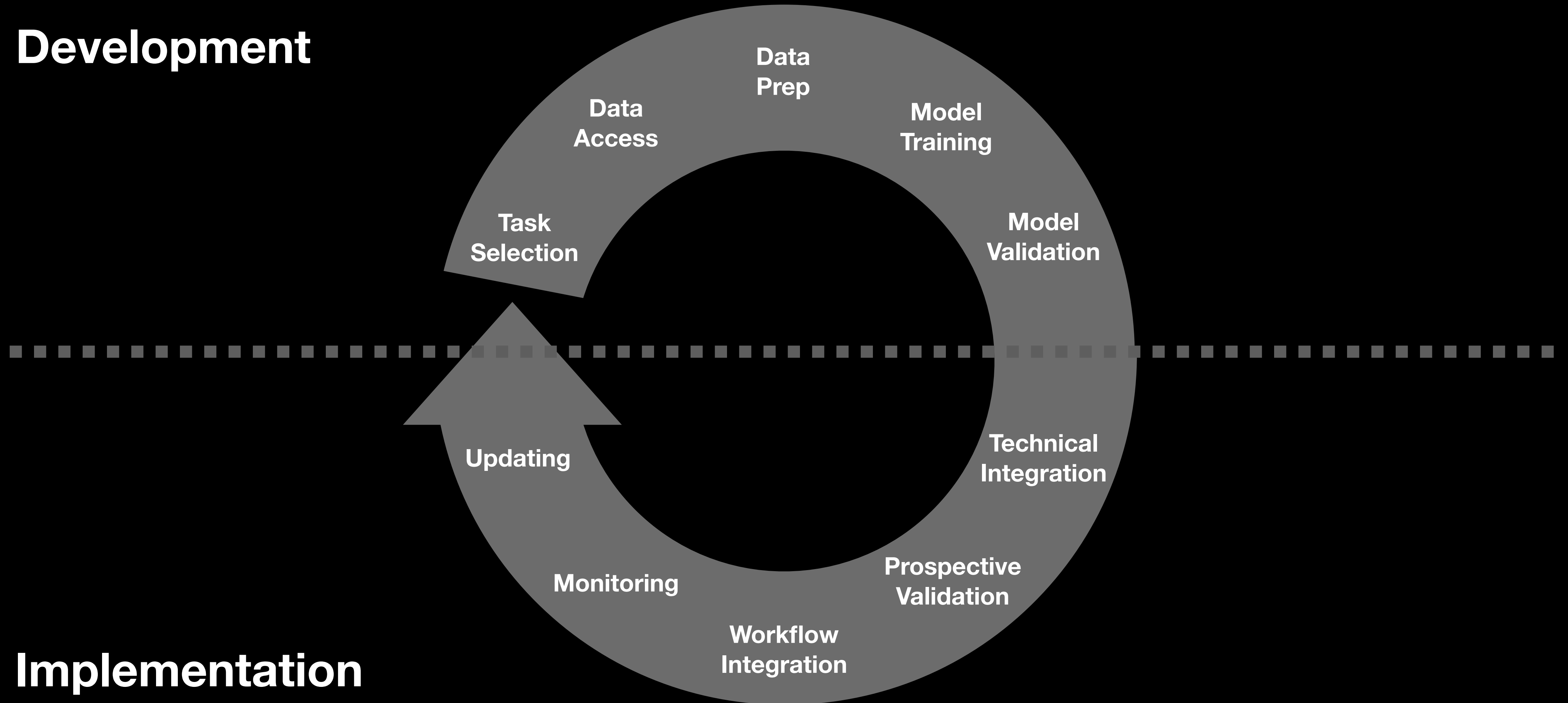
Integration into care

## Implementation

Thesis: by focusing on the longitudinal nature of healthcare, we can improve the performance of predictive models during development & implementation.

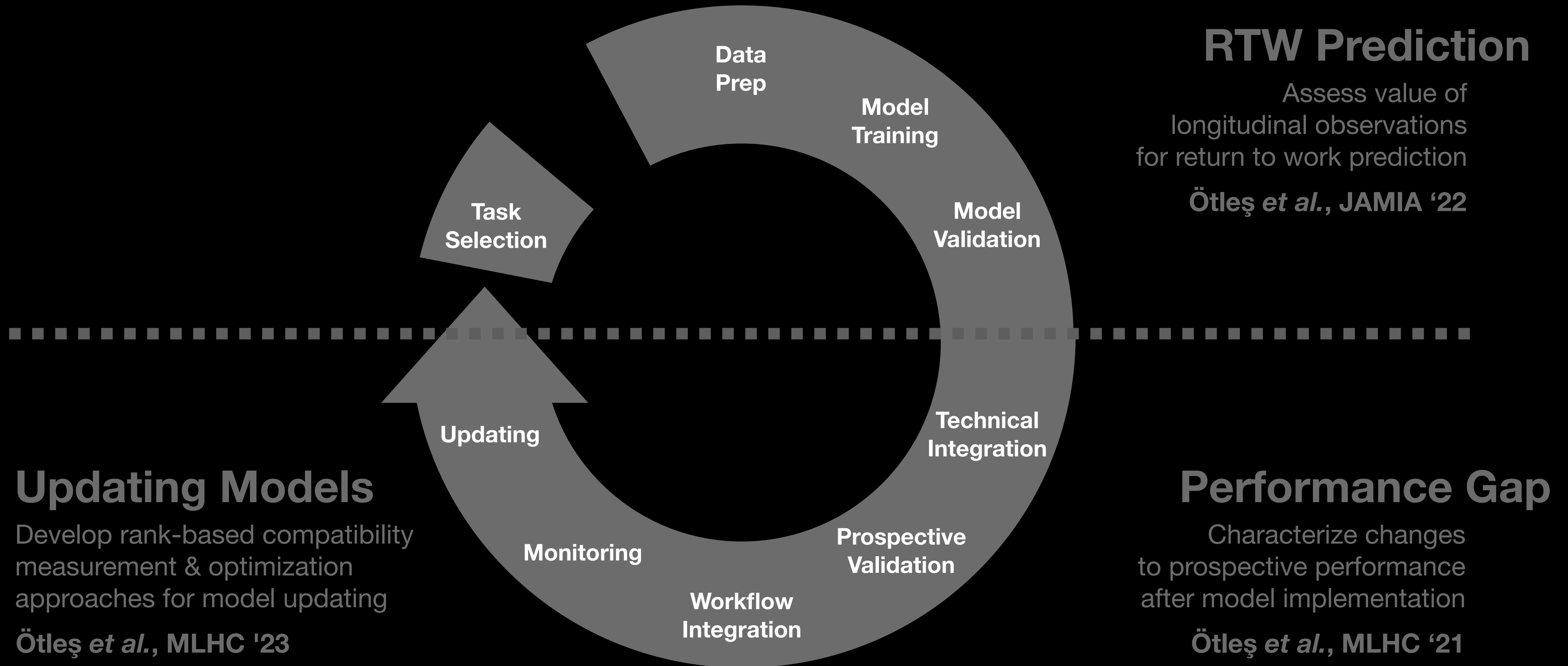
# The ML Lifecycle

**Development**



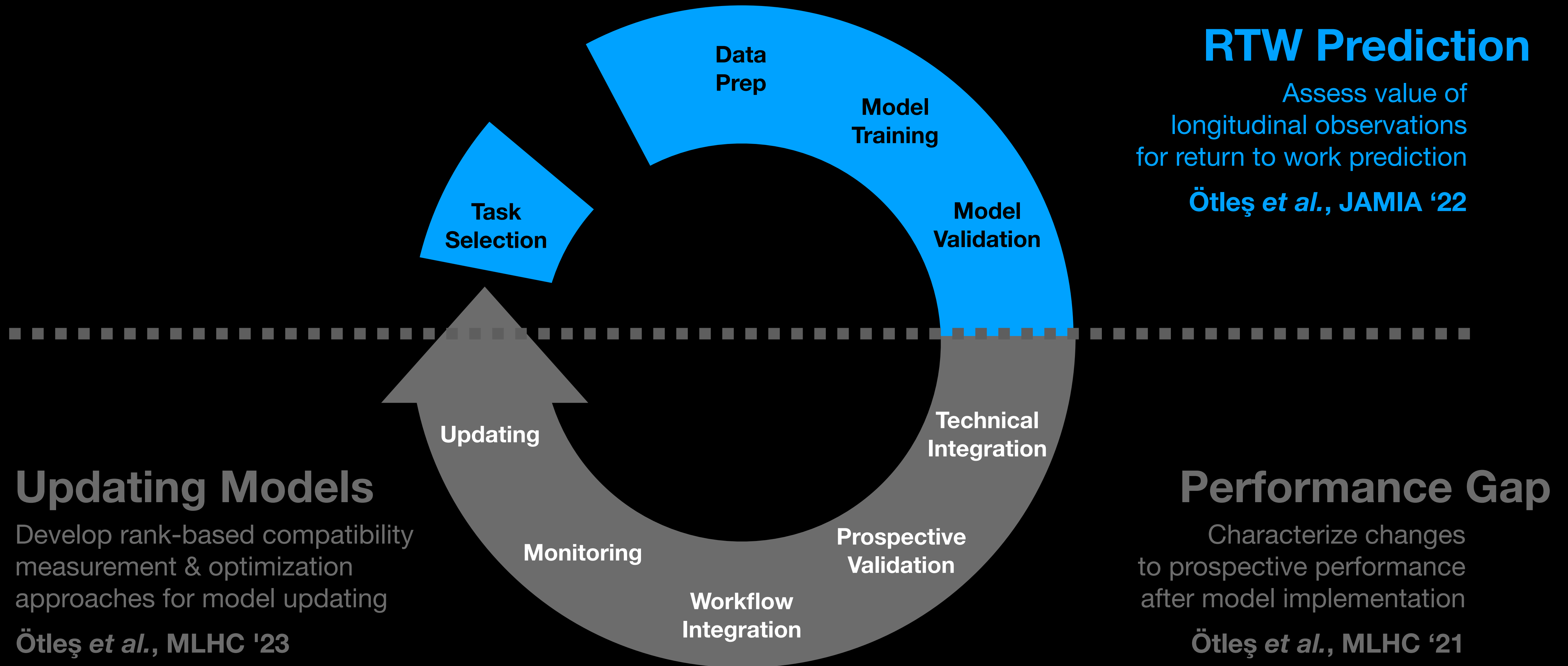
**Implementation**

# The ML Lifecycle

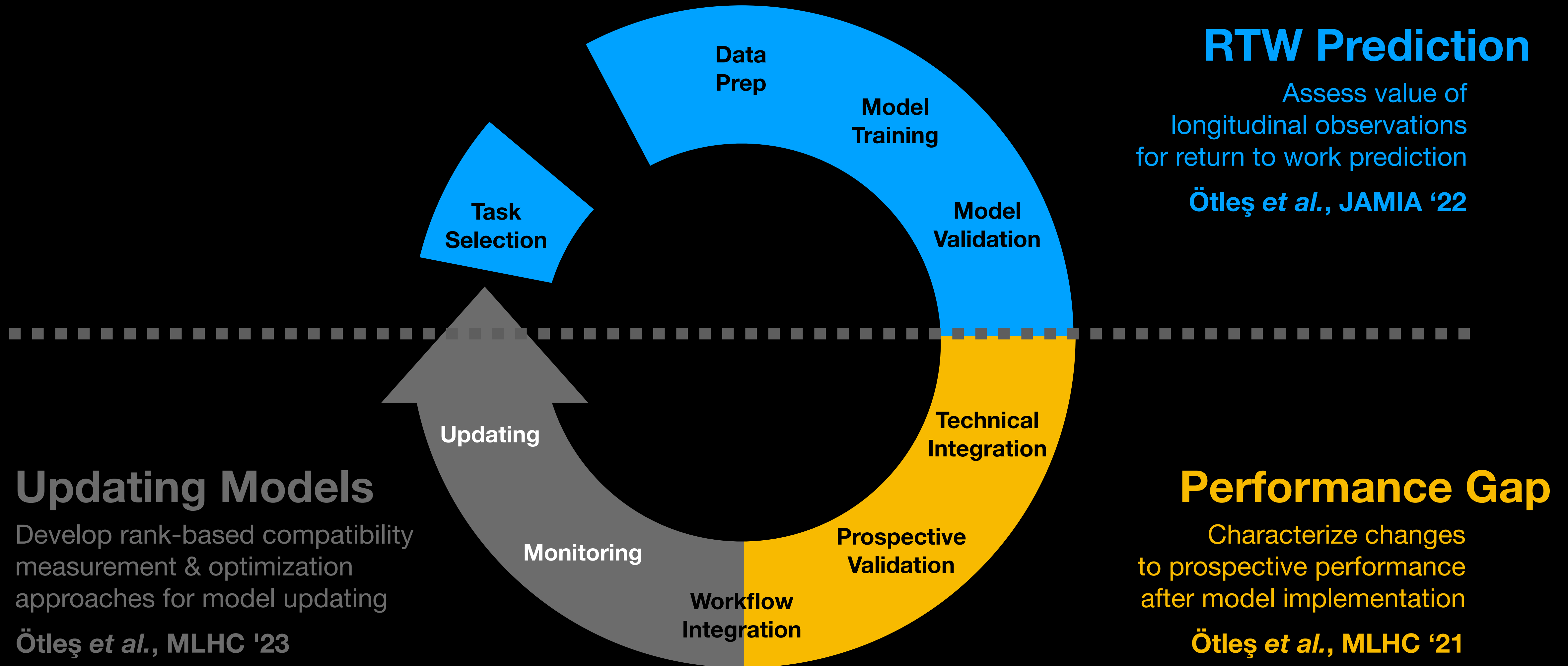




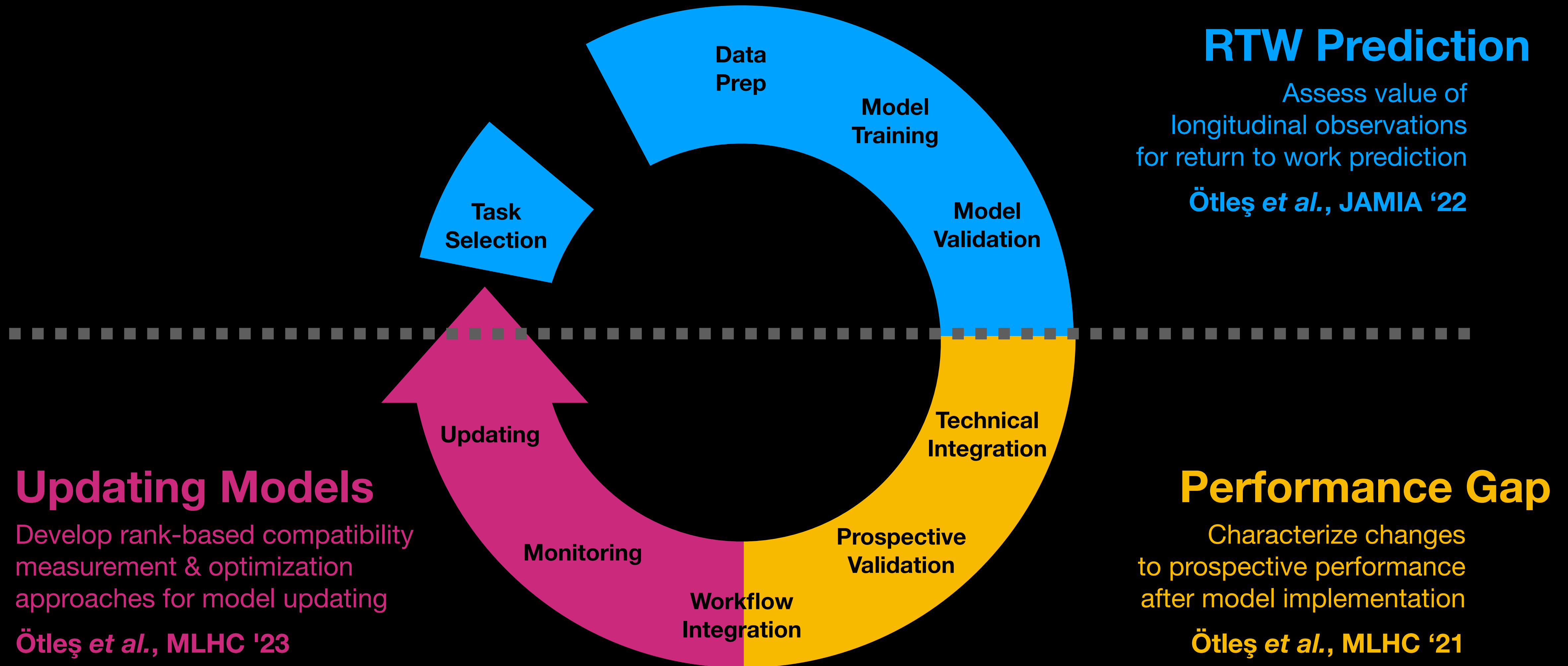
# The ML Lifecycle



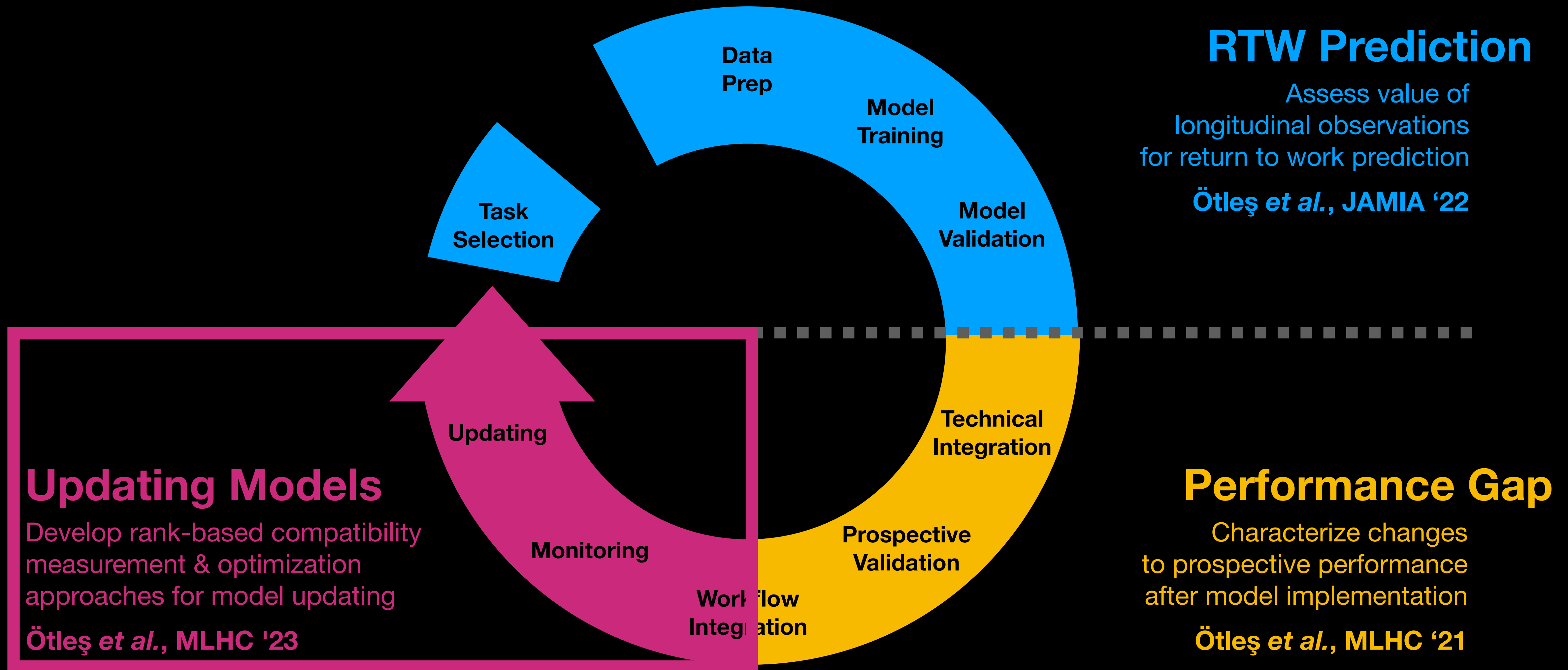
# The ML Lifecycle



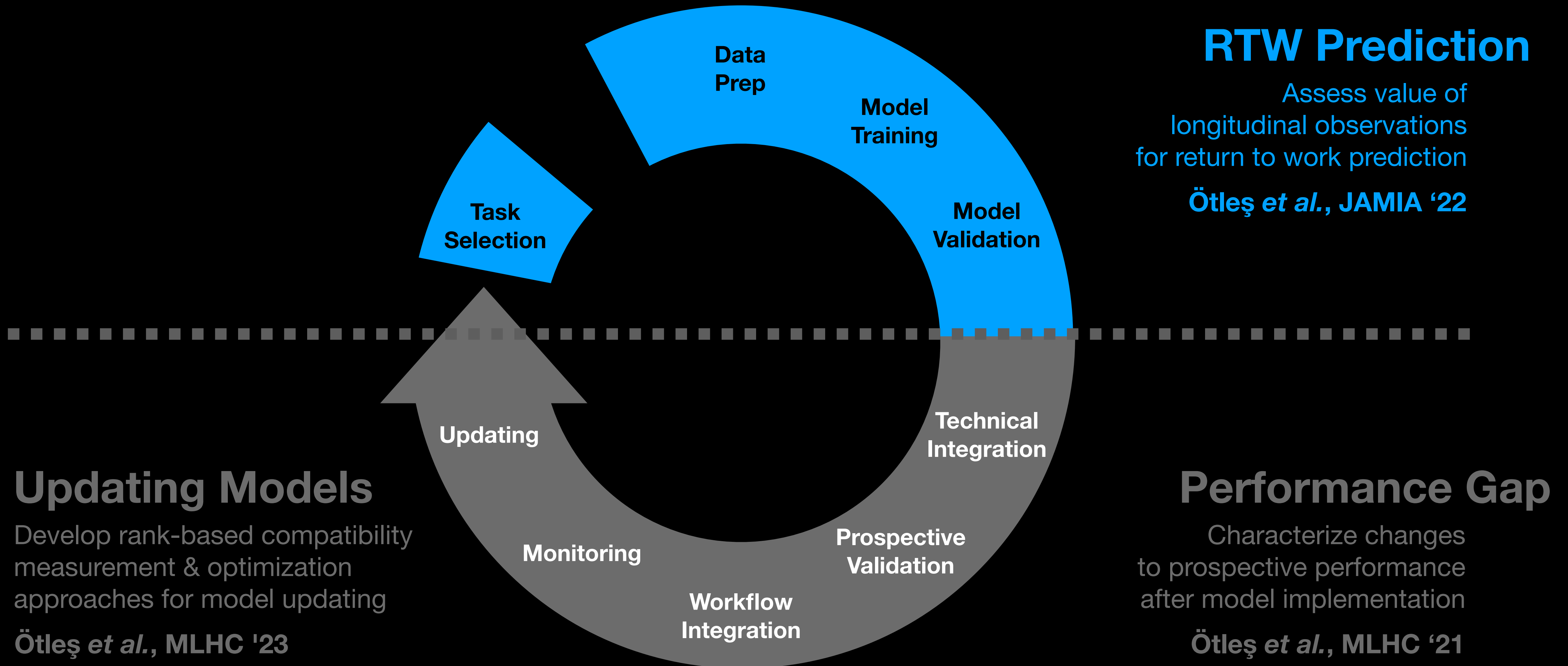
# The ML Lifecycle



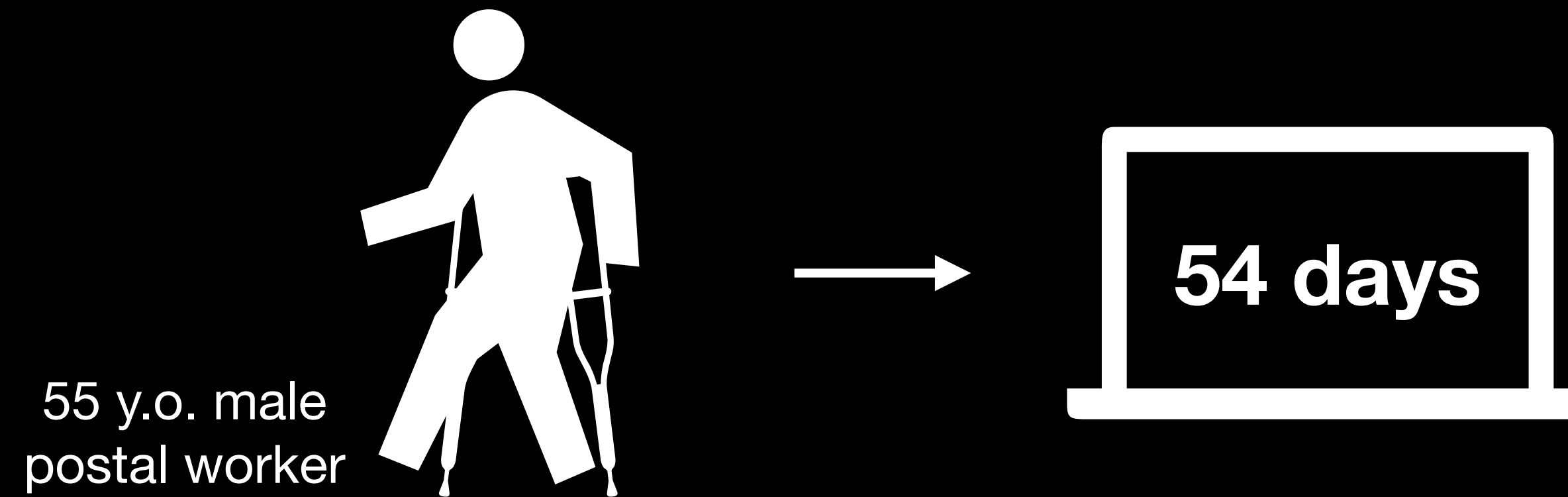
# What are we going to talk about?



# Let's dive in.



# Existing return to work models ignore longitudinal observations.

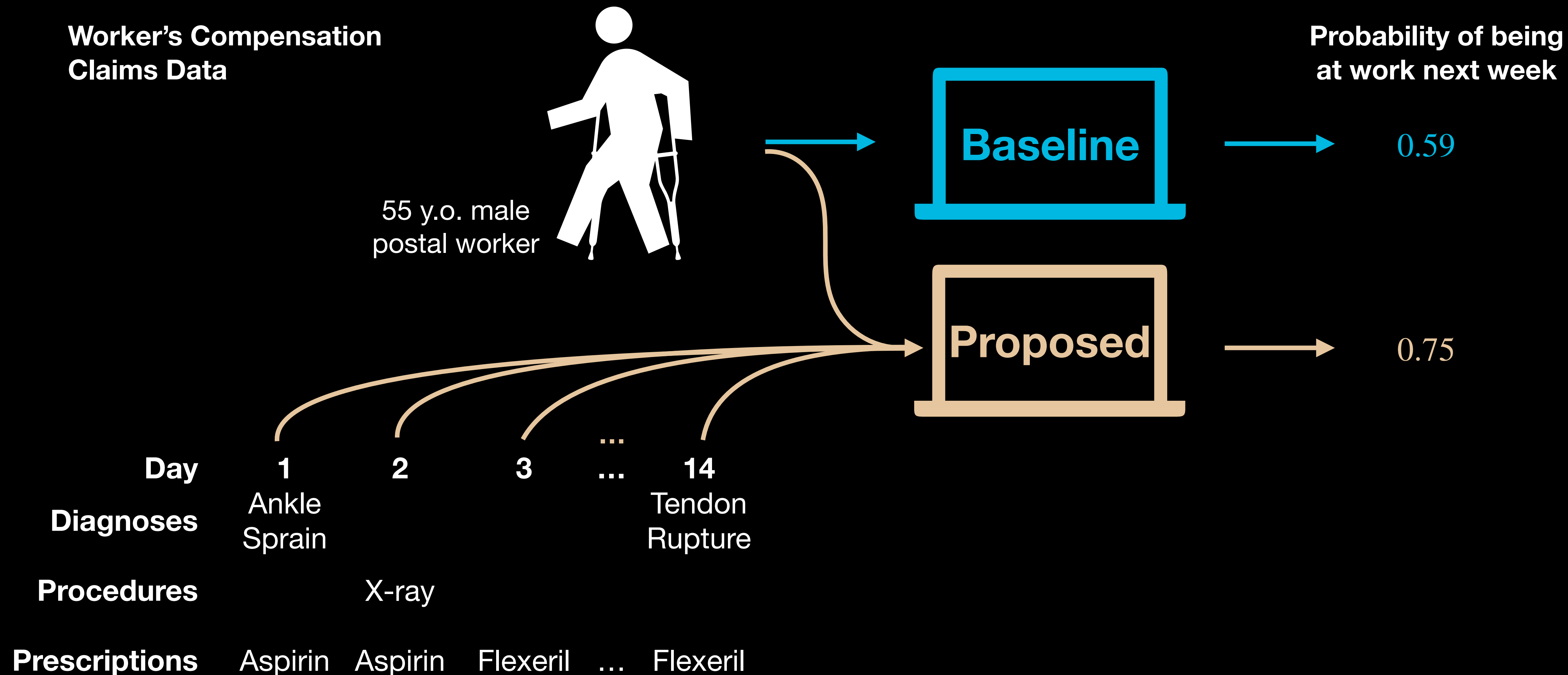


***Q: What is the value of longitudinal observations in return to work prediction?***

More specifically: when predicting future work-status do we observe a performance improvement when using longitudinal observations collected beyond the time of injury?

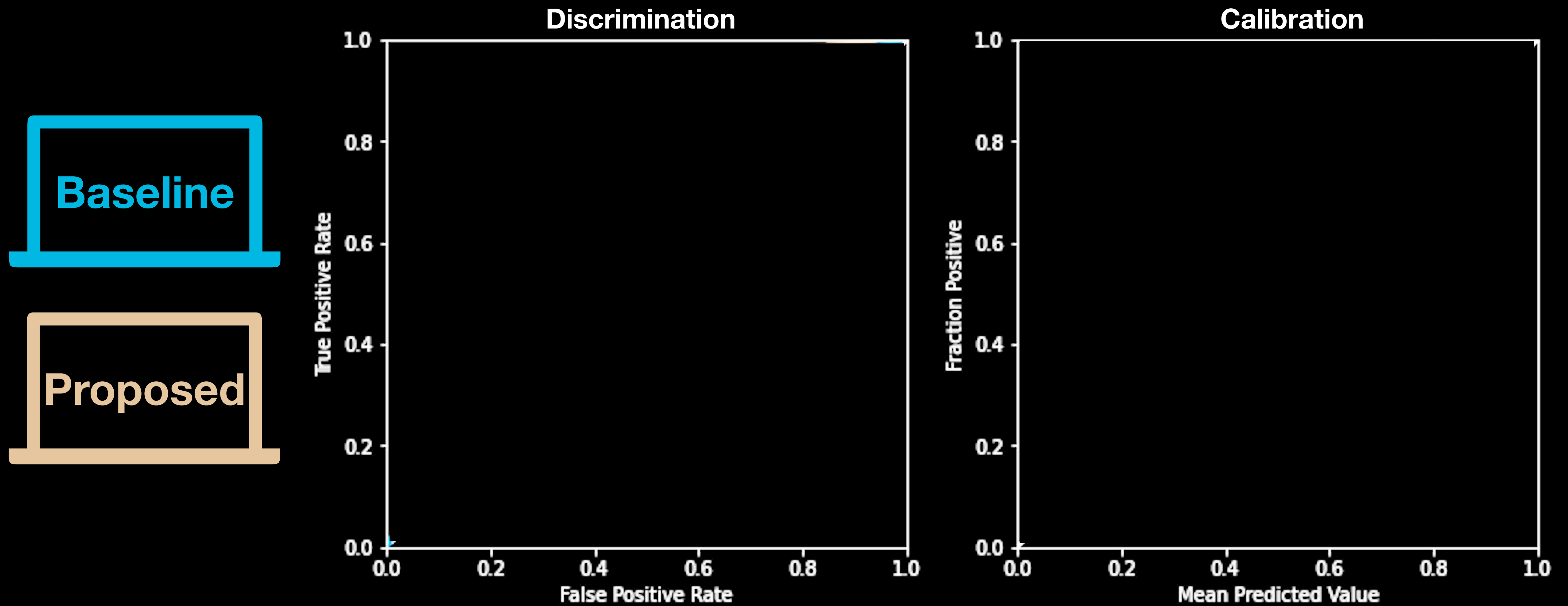
Hypothesis: Yes. Longitudinal observations improve predictions in other healthcare task.

# Experimental Setup

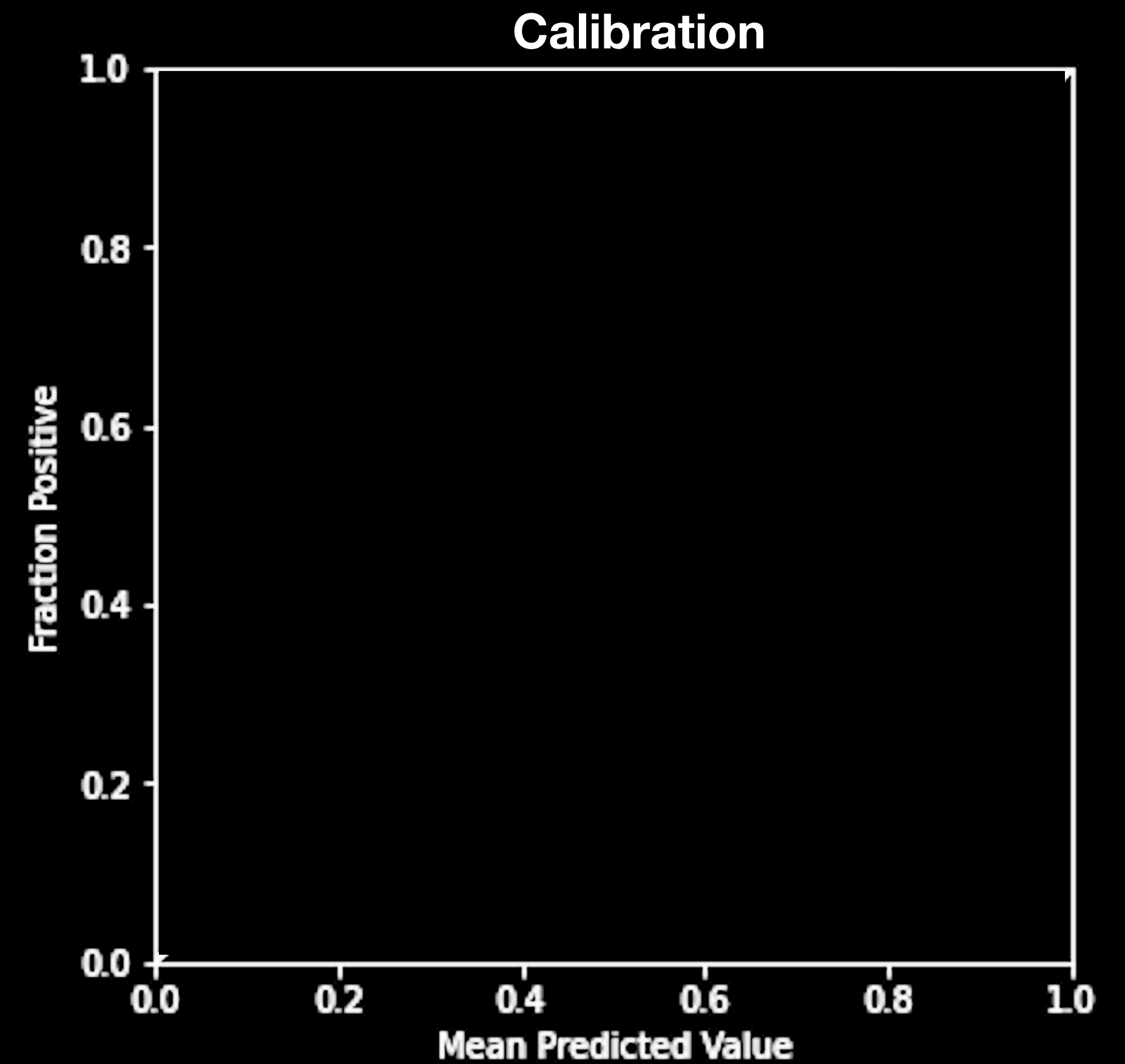
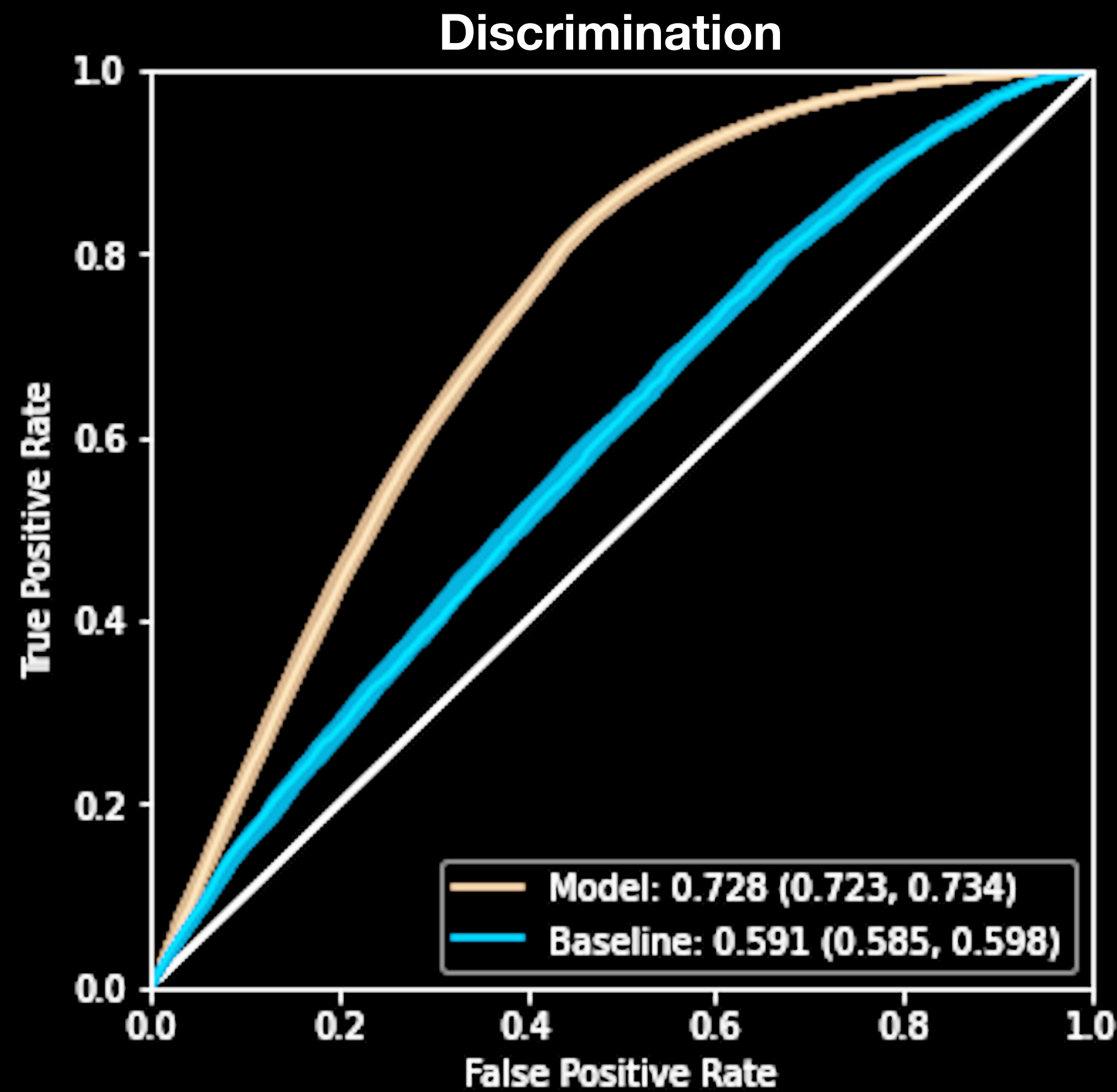
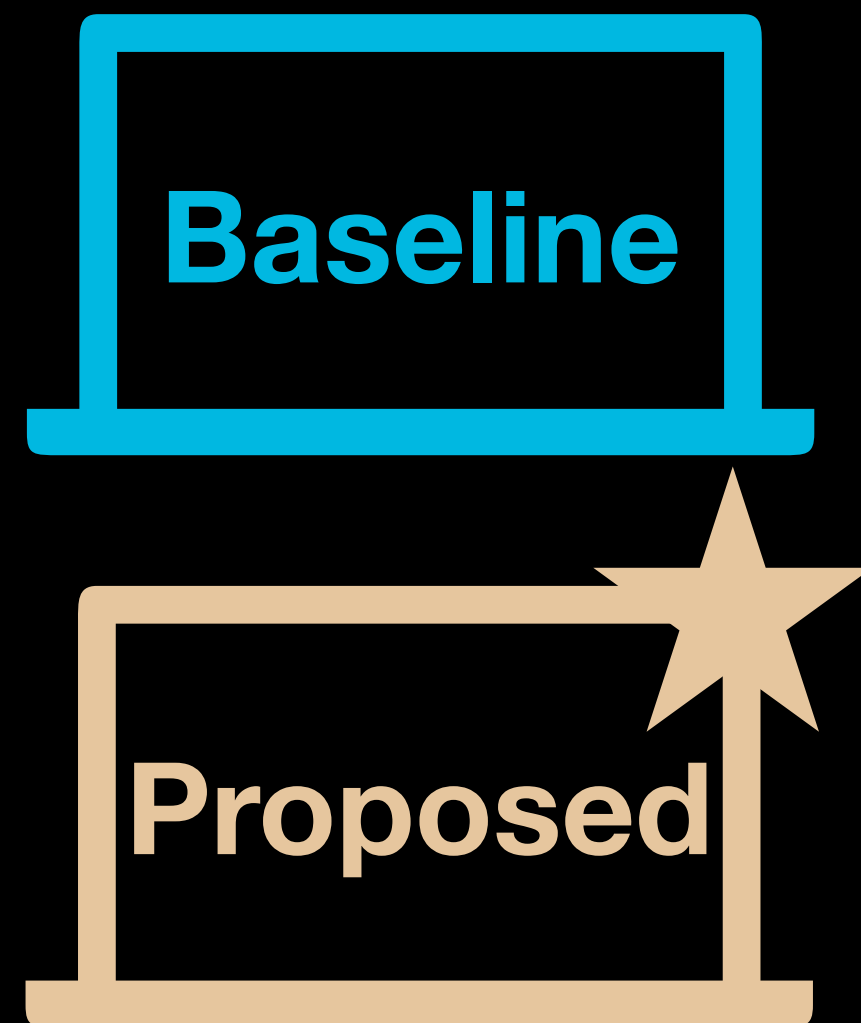




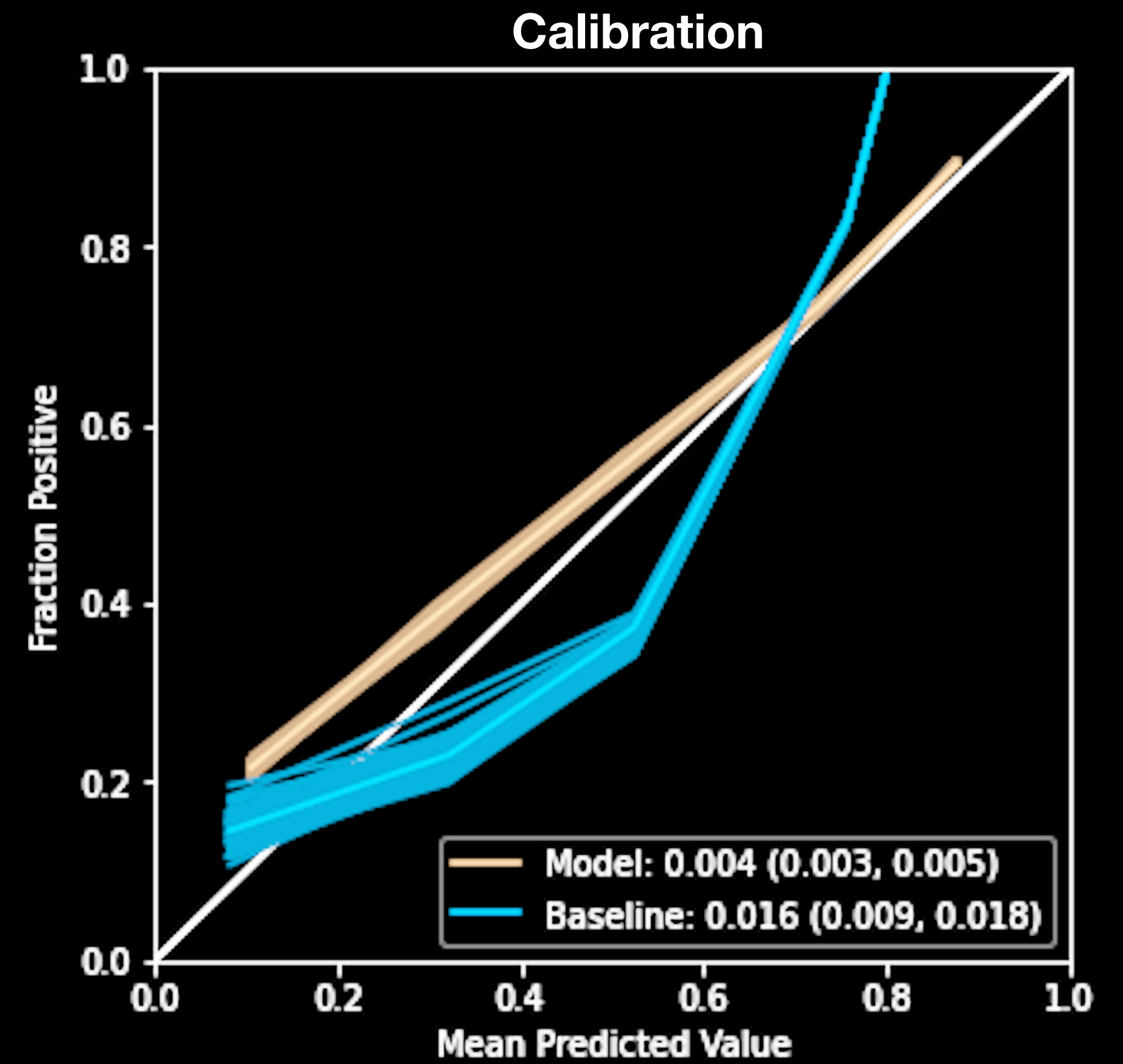
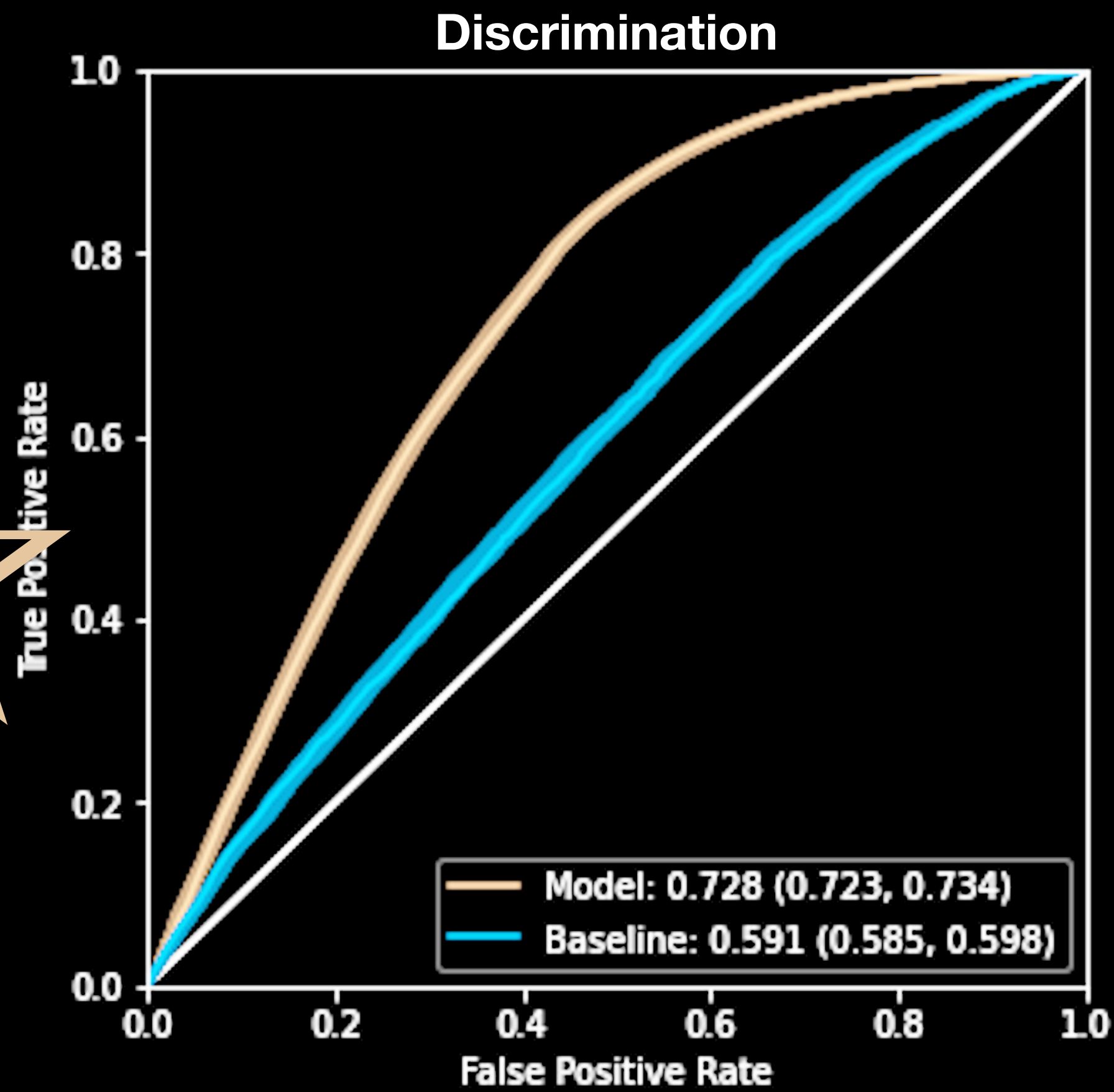
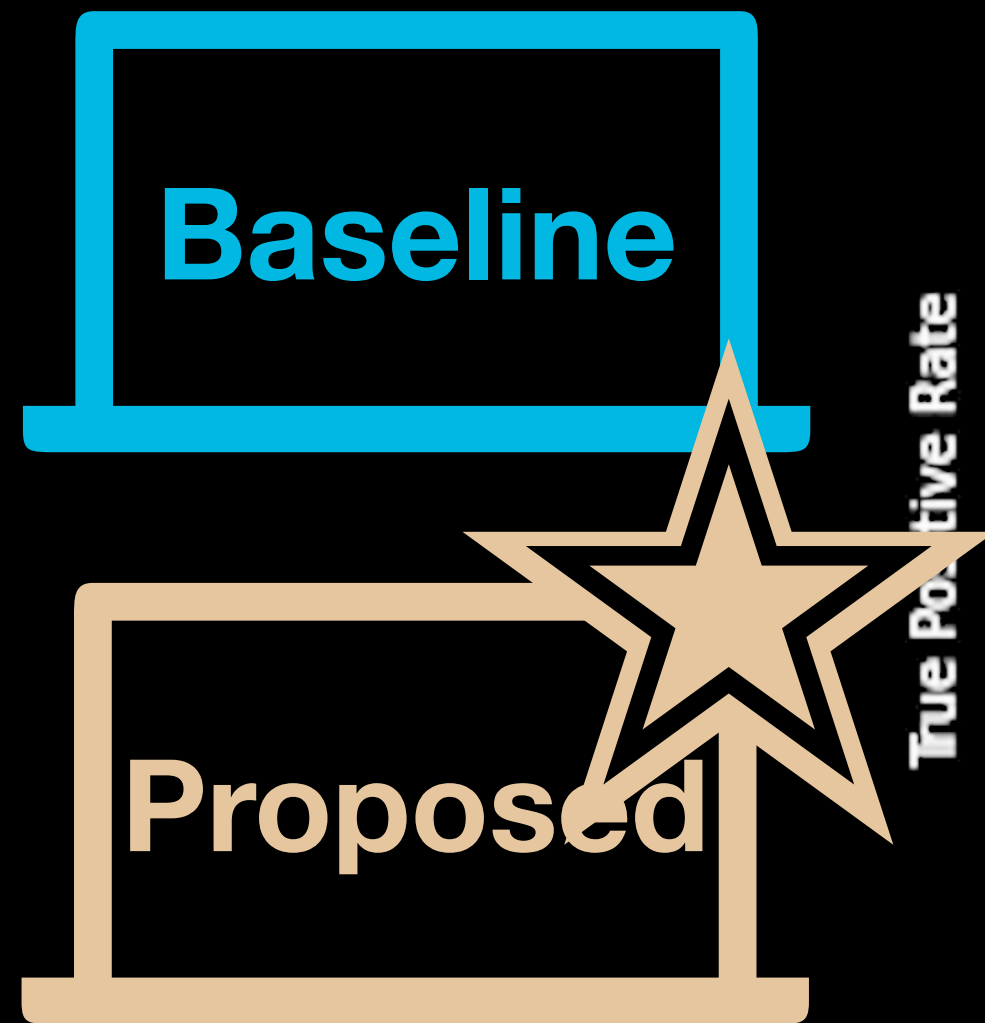
# Results



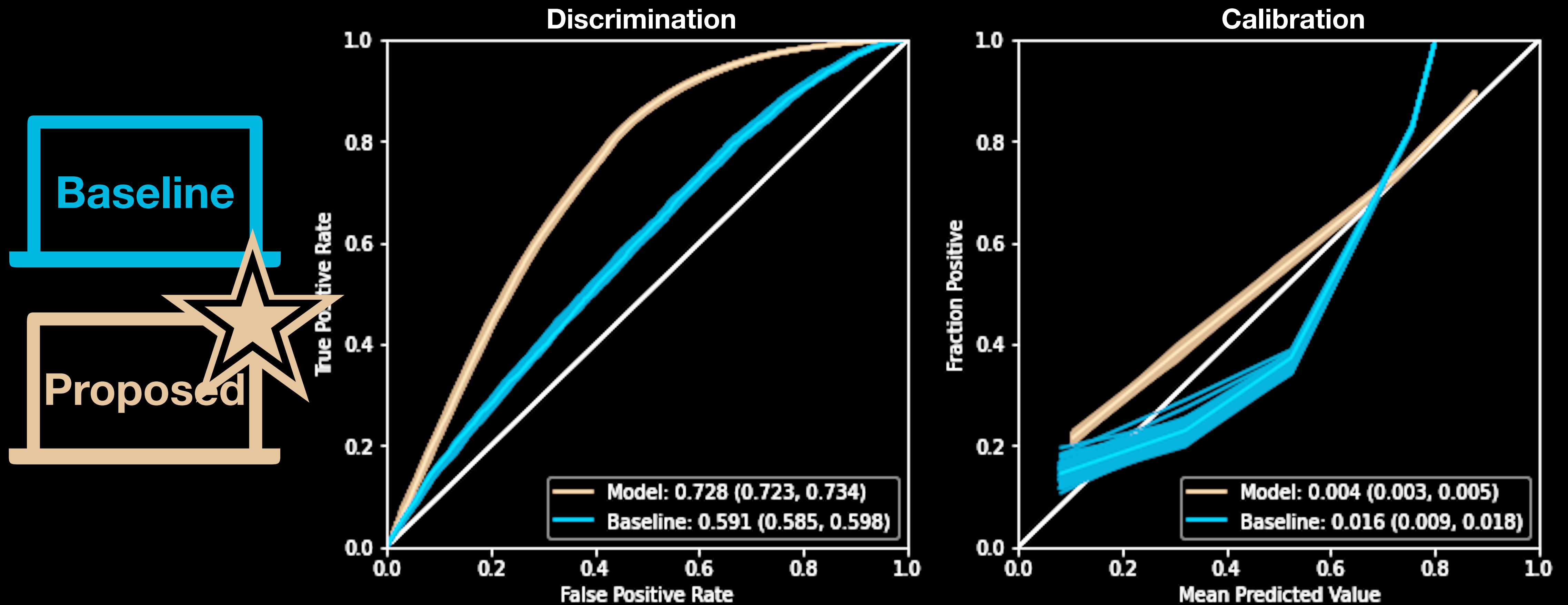
# Results



# Results

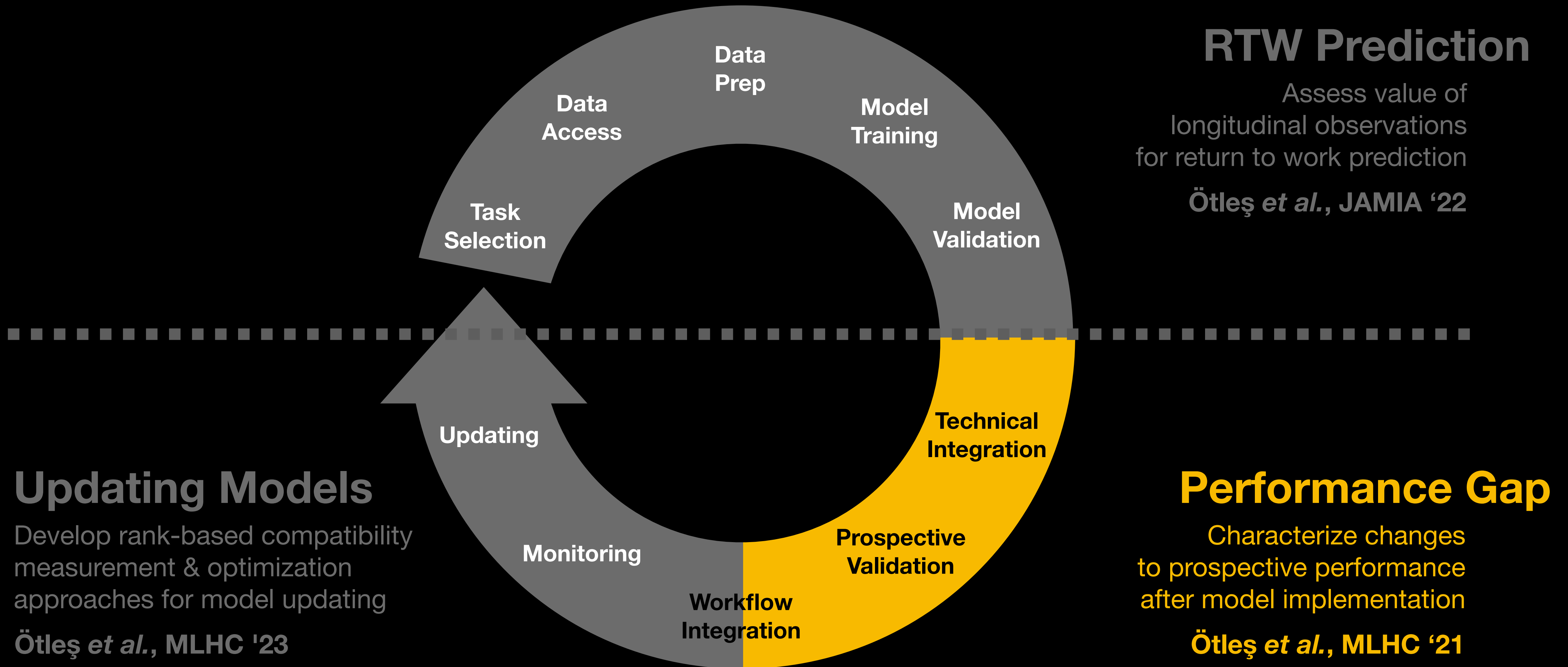


# Main Takeaway from RTW Prediction.

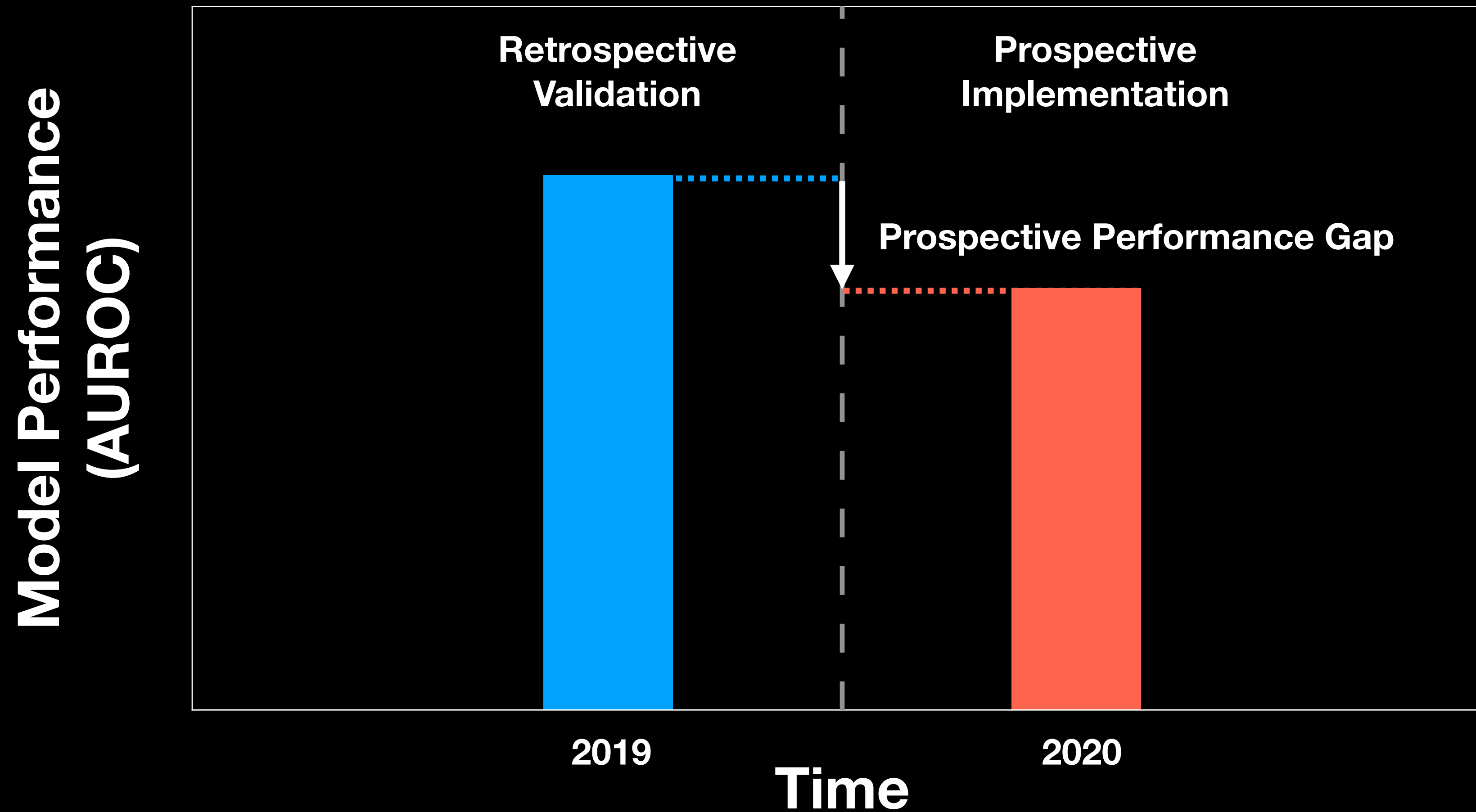


Return to work prediction benefits from longitudinal observations.

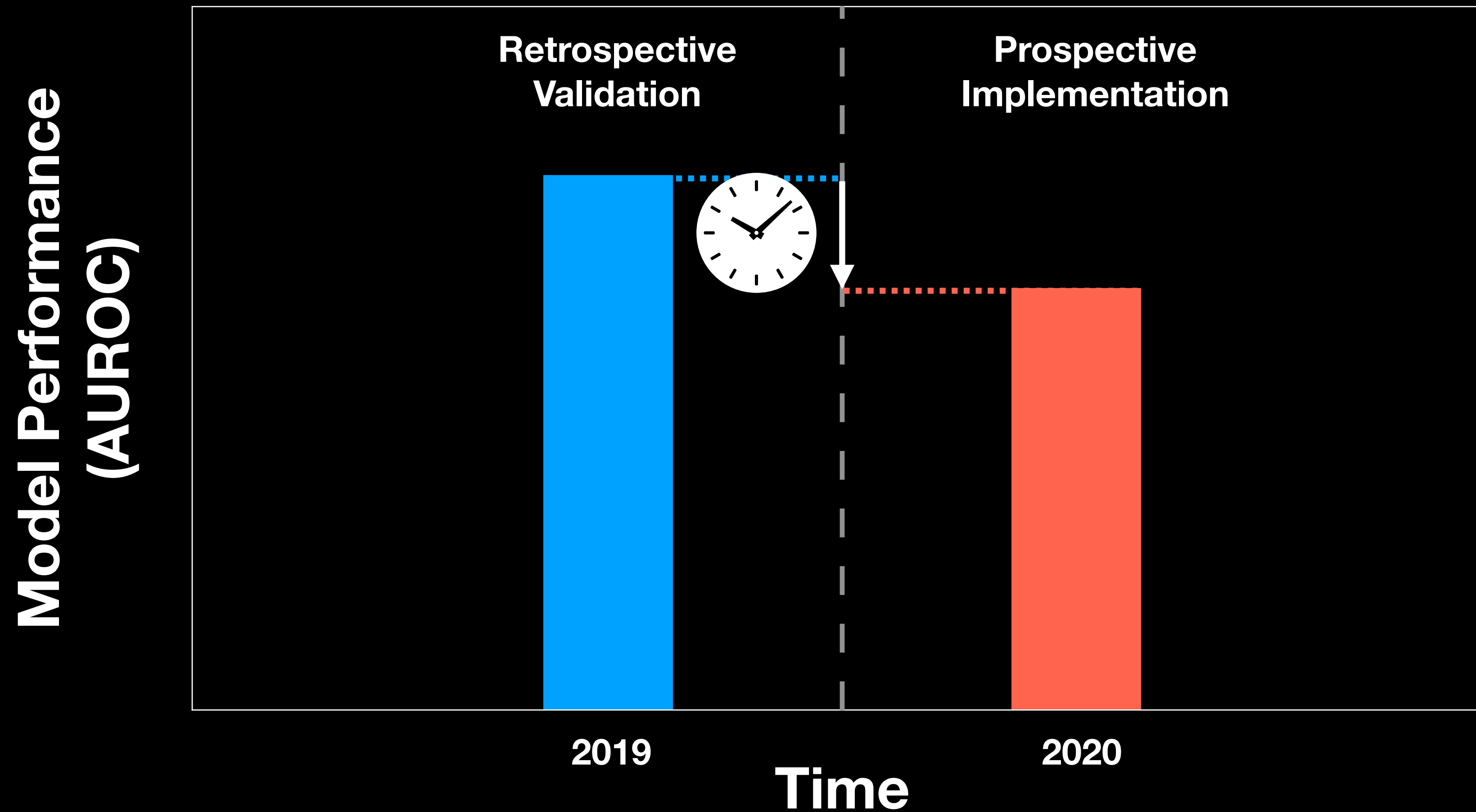
# On to the next one.



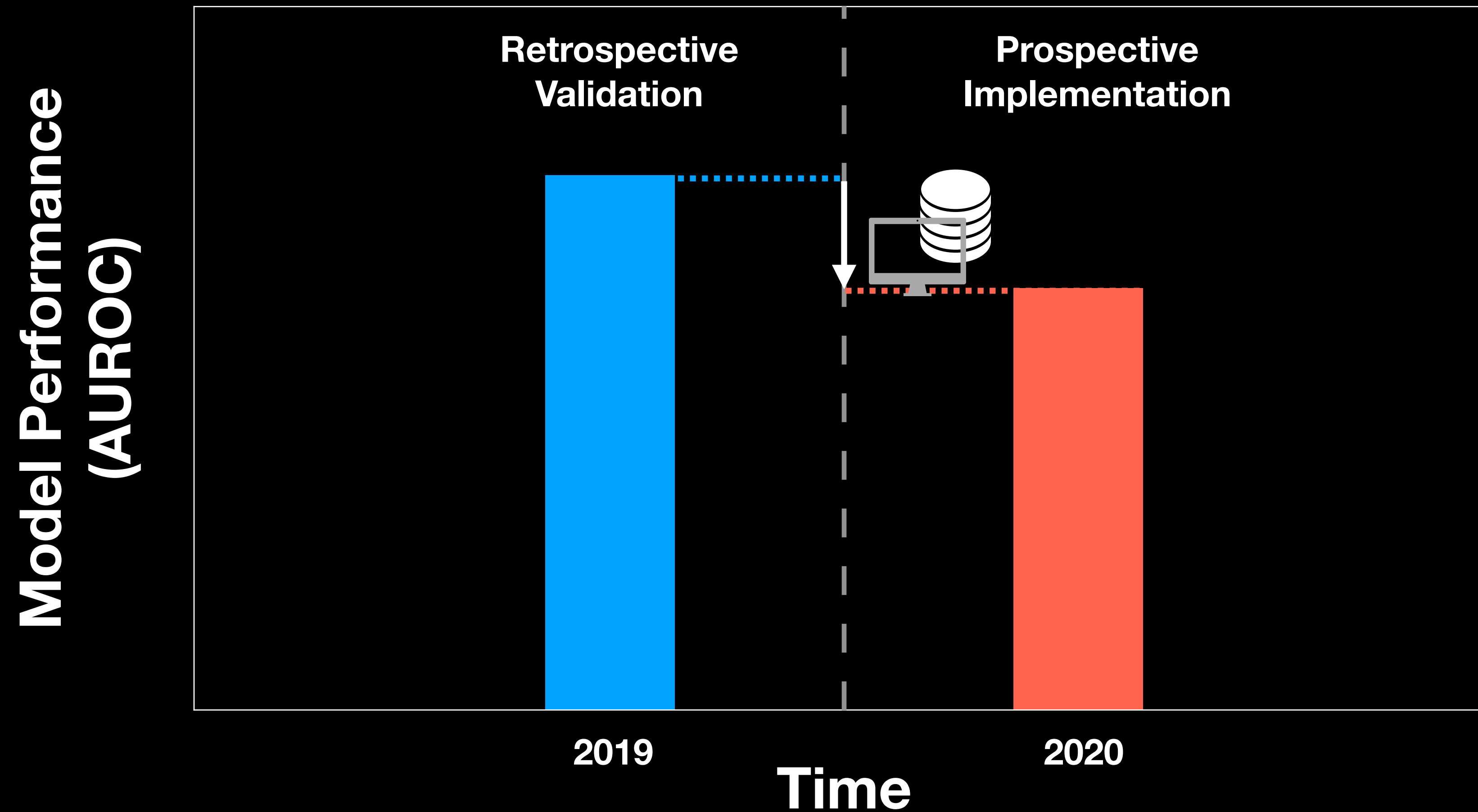
# Model performance may degrade after implementation.



This degradation is often attributed to changes in populations & practice that occurs over time.

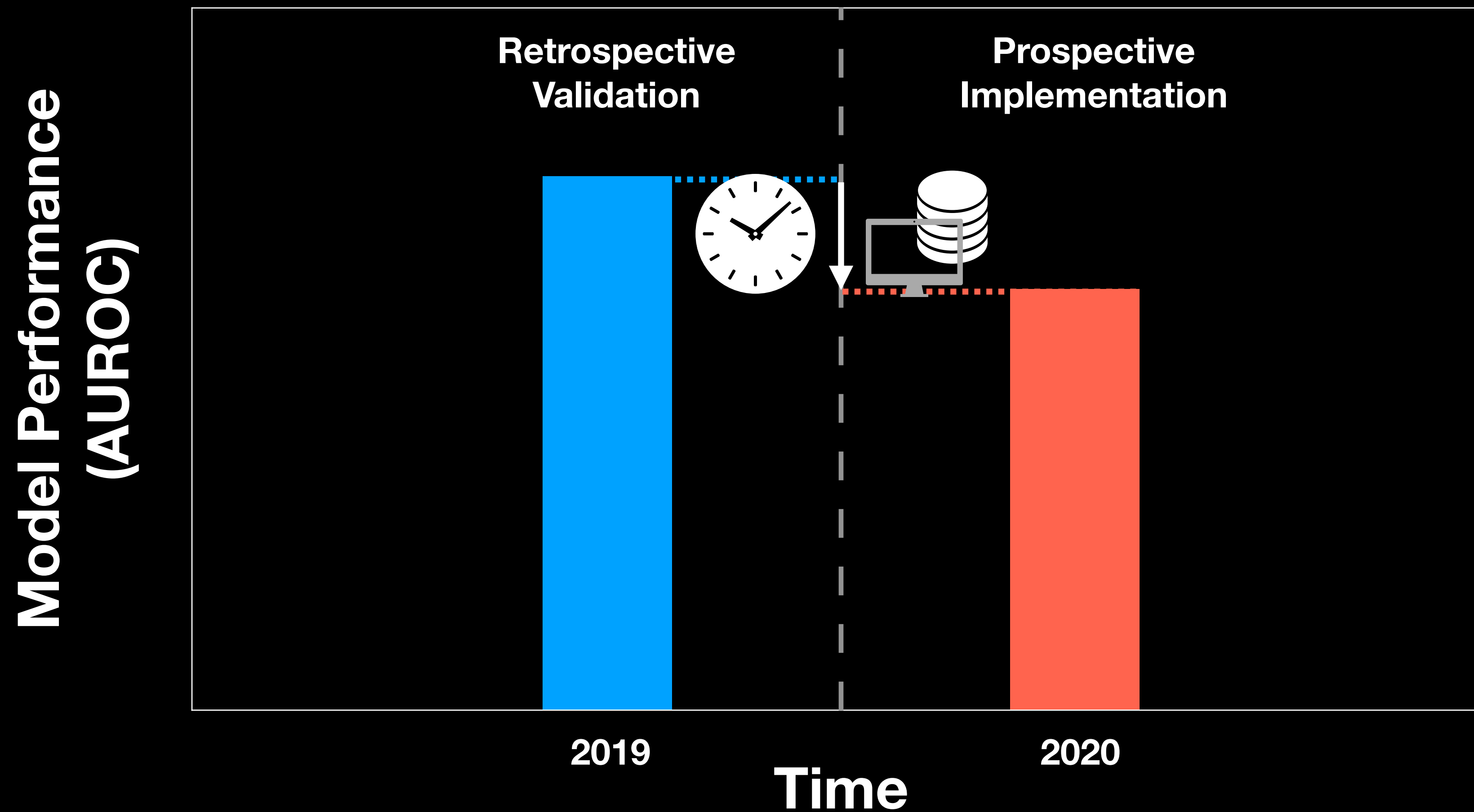


However, changes in IT infrastructure may also affect the prospective performance gap.





# Degradation due to temporal & infrastructure shift.

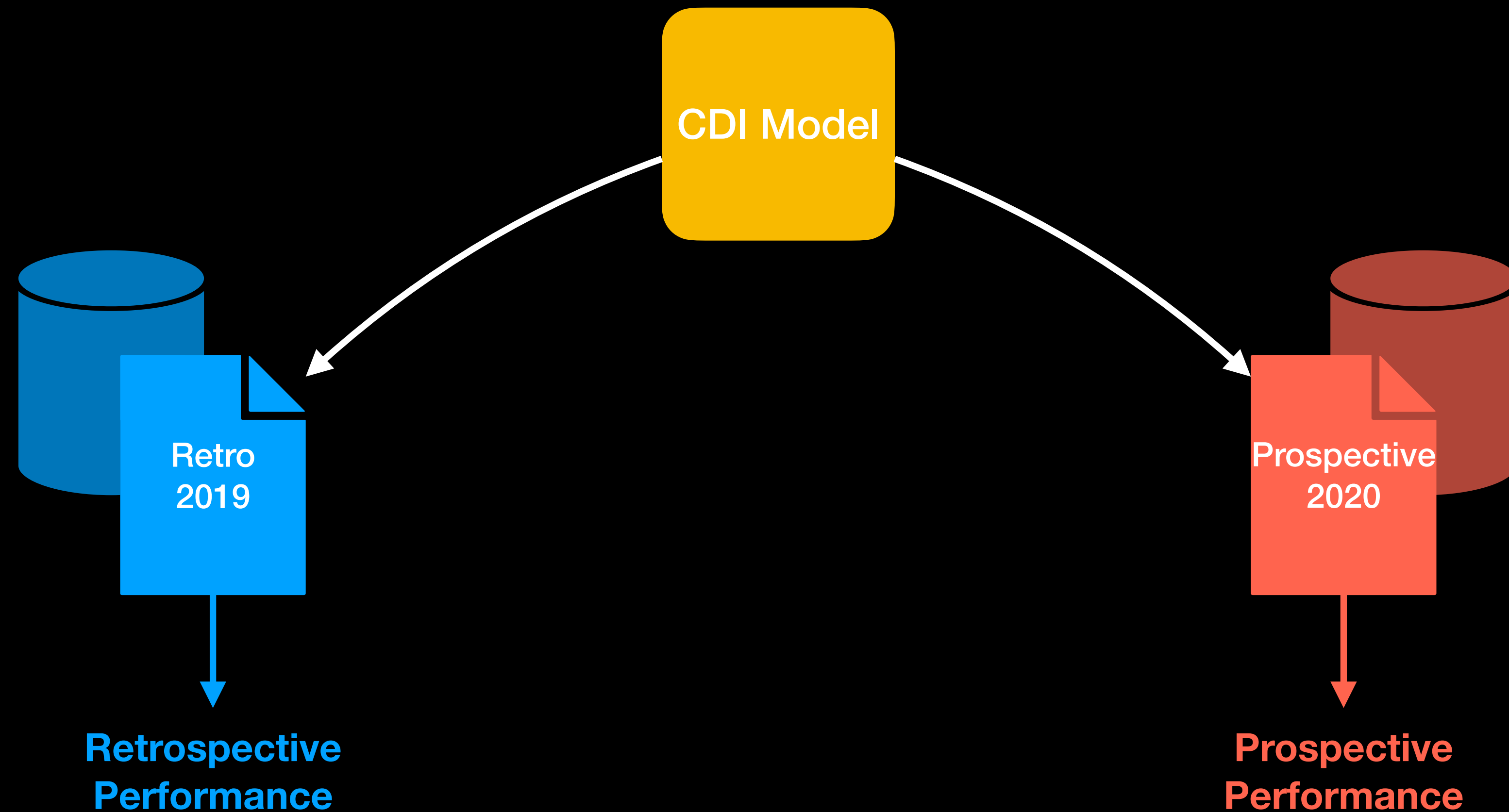


***Q: What are the causes of the prospective performance gap?***

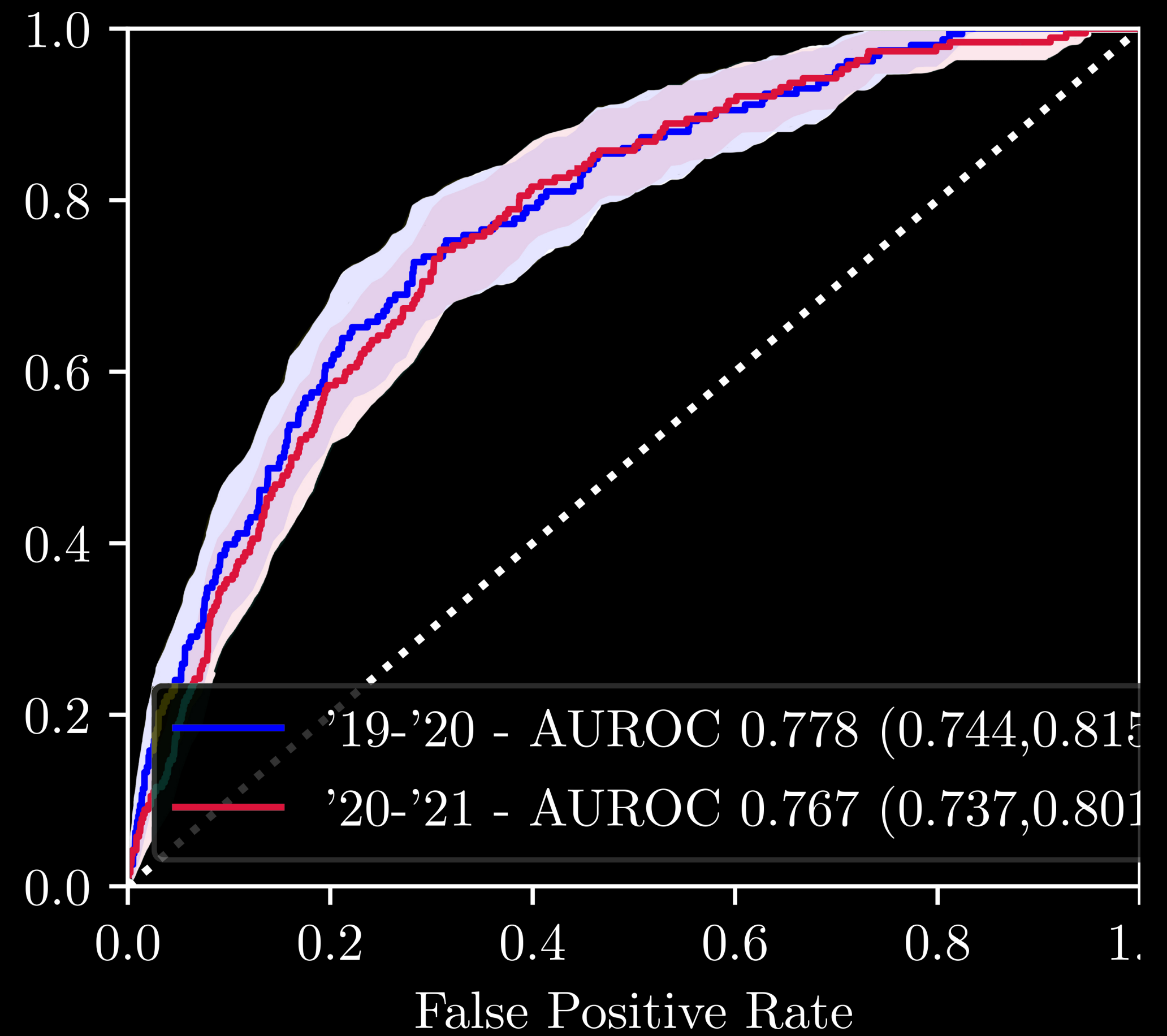
More specifically: what are the temporal and infrastructure shift contributions to the prospective performance gap of a prospectively implemented model?

Hypothesis: Infrastructure will be an important component of the prospective performance gap.

# Experimental Setup

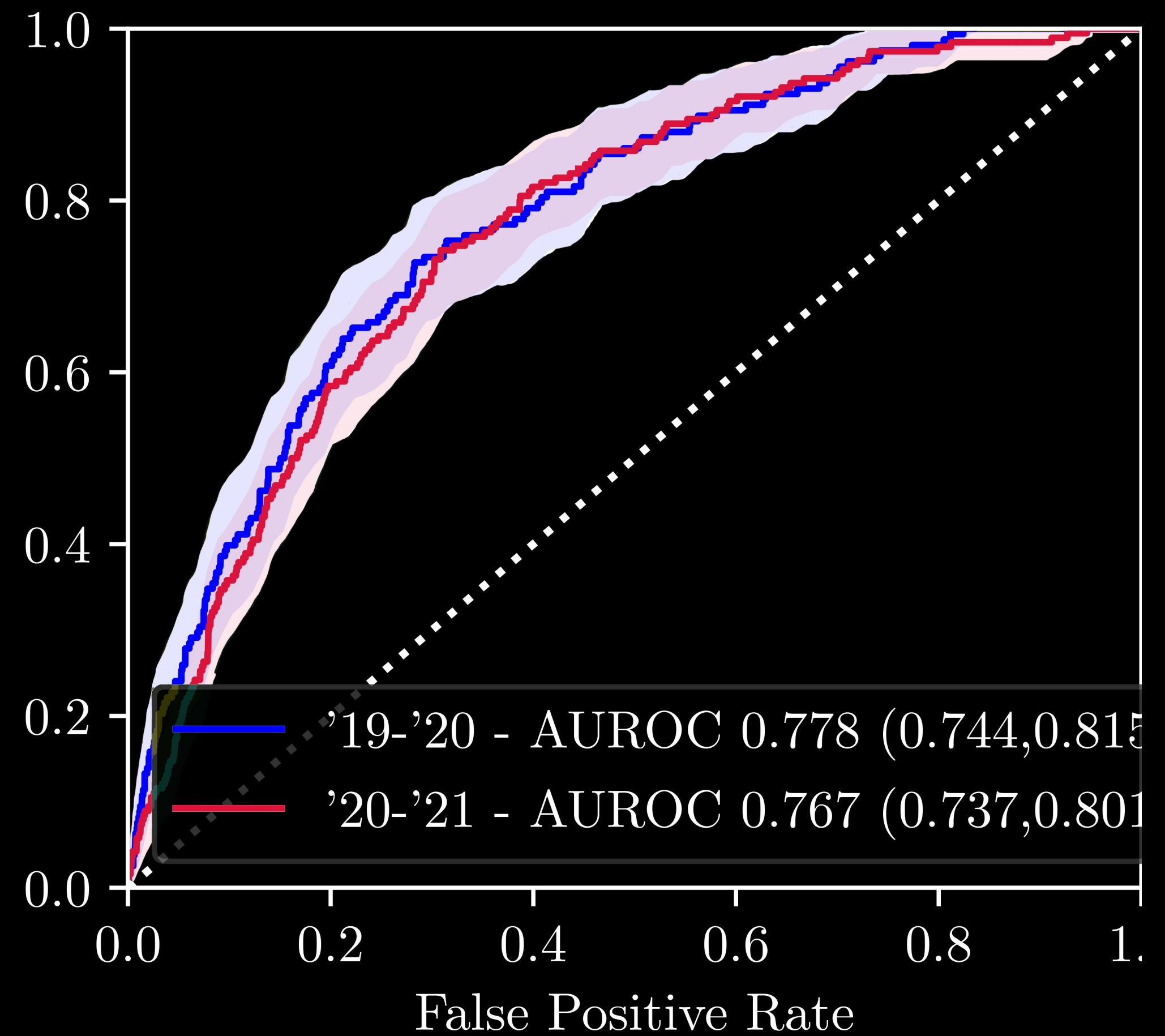


# Results



# Results

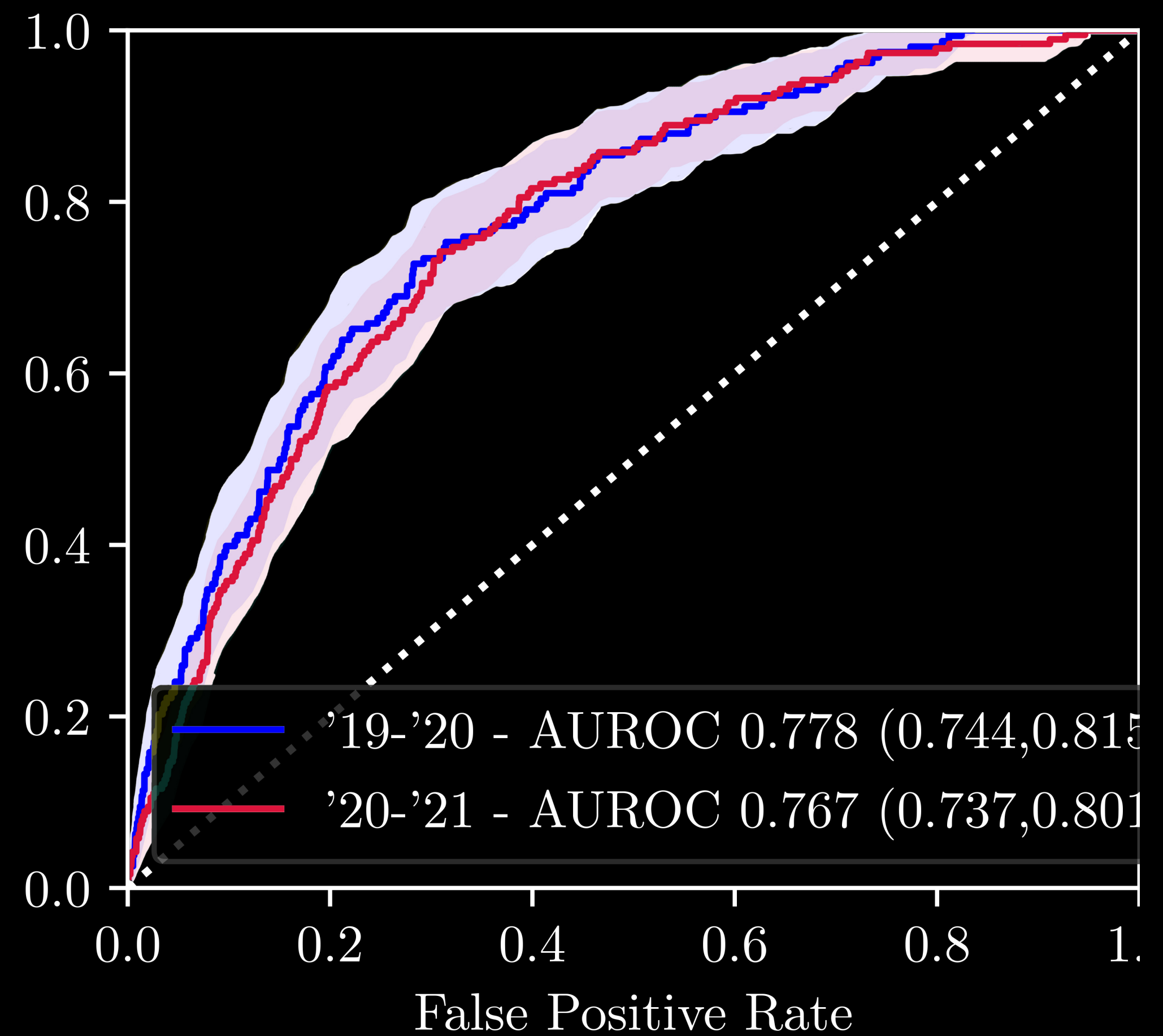
Performance Gap: **0.011**



# Results

Performance Gap: **0.011**

Temporal Shift: **-0.005**

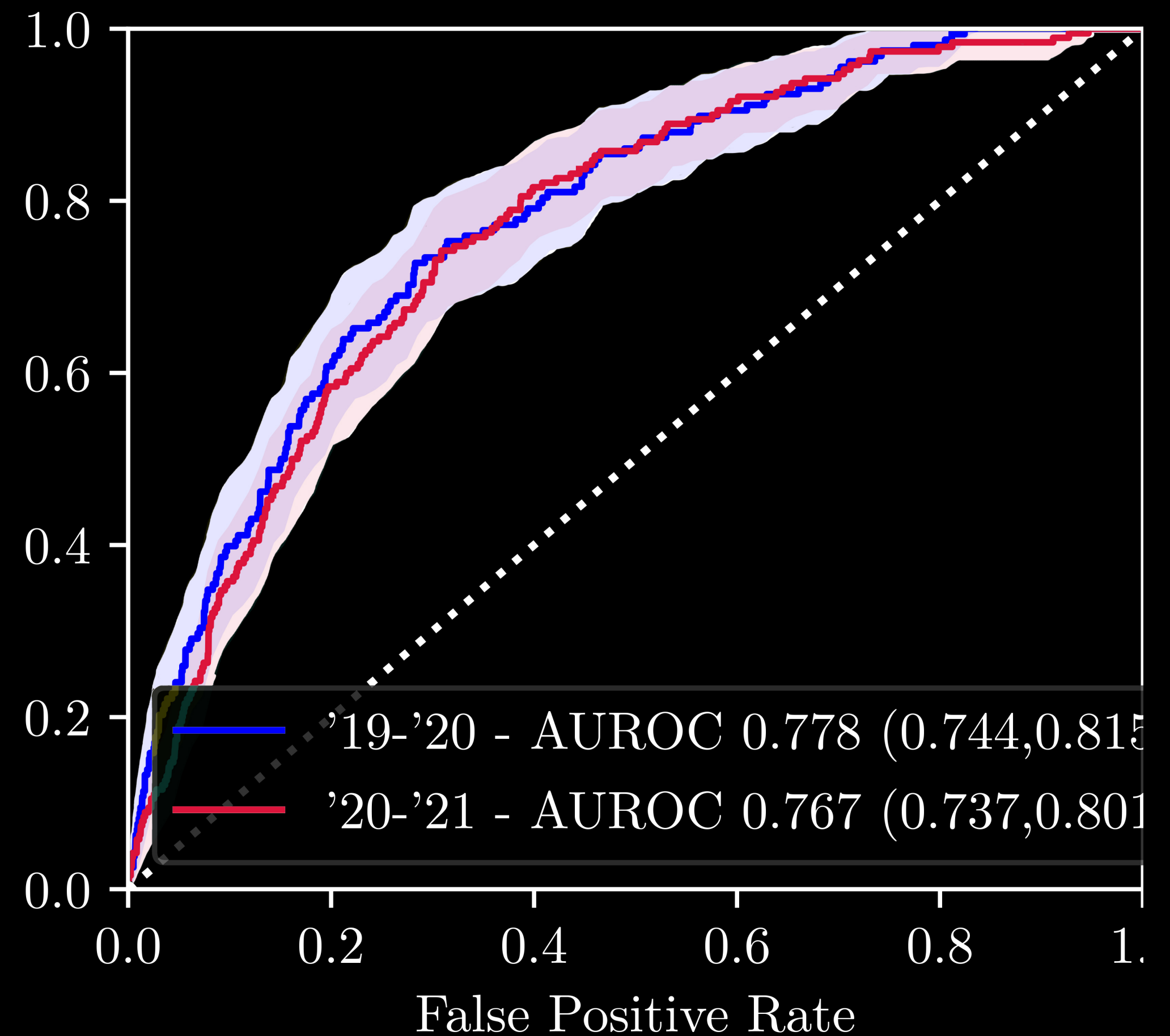


# Results

Performance Gap: **0.011**

Temporal Shift: **-0.005**

Infrastructure Shift: **0.016**



# Infrastructure Matters.

Observed prospective performance gap primarily driven by infrastructure shift:

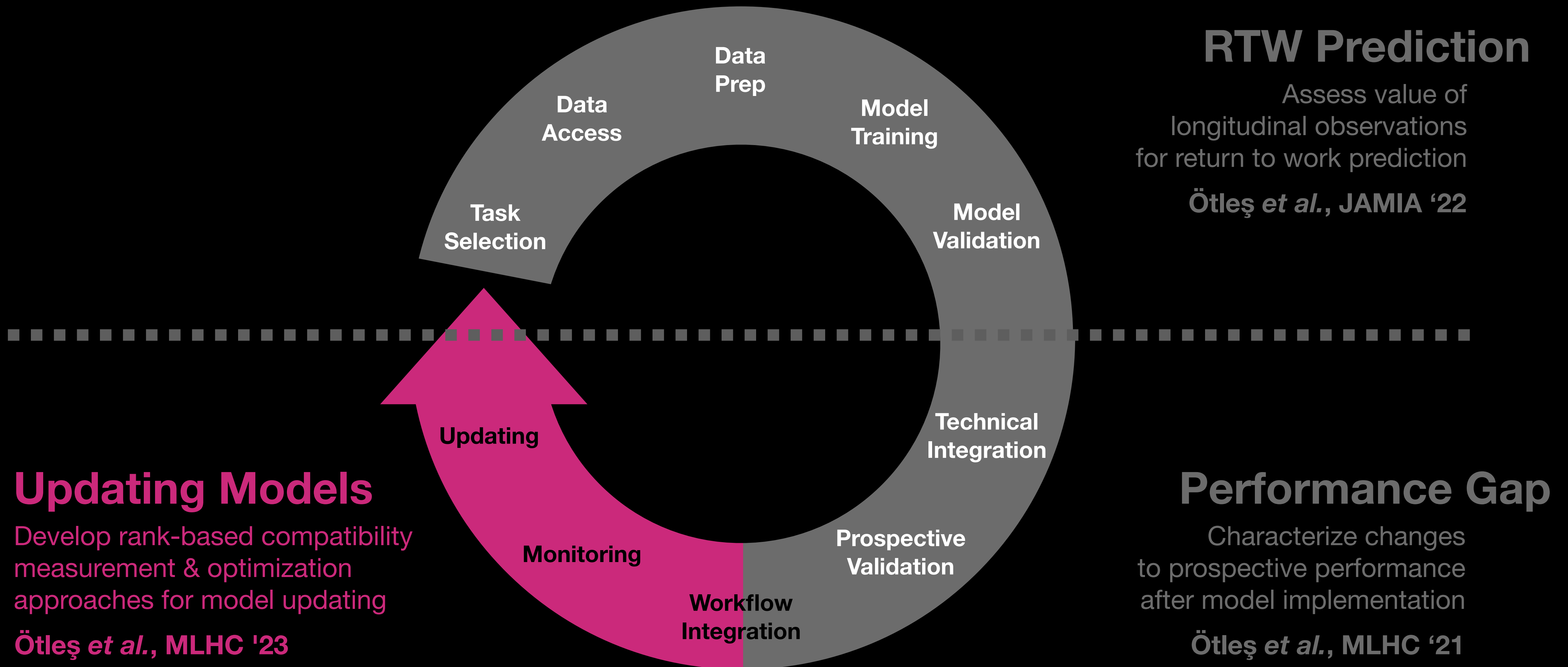
source & processing differ, requires infrastructure updates

timing of access differs, *when* you access the data affects *what* you access

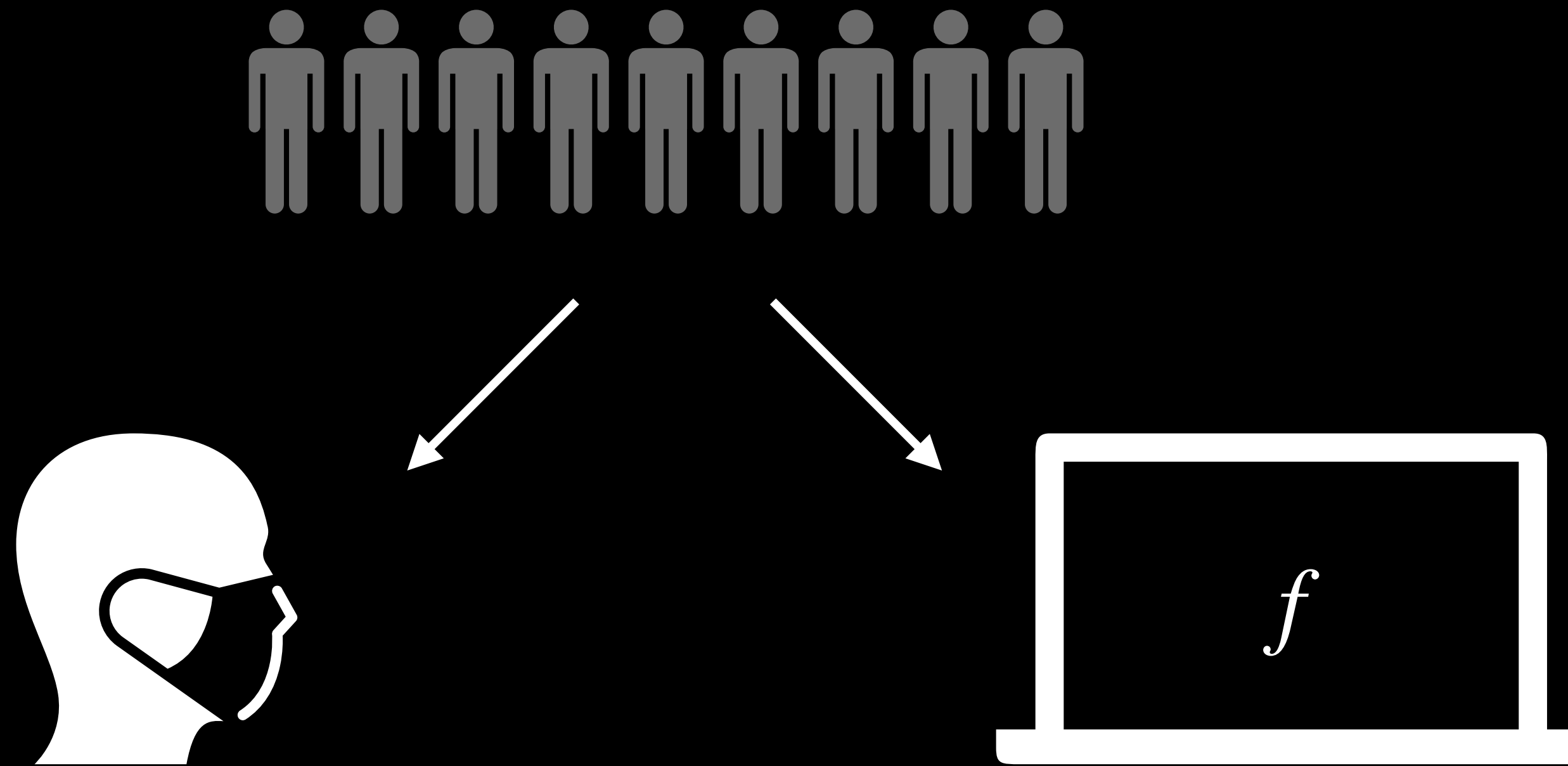
**Infrastructure shift can be mitigated, but first it must be recognized!**



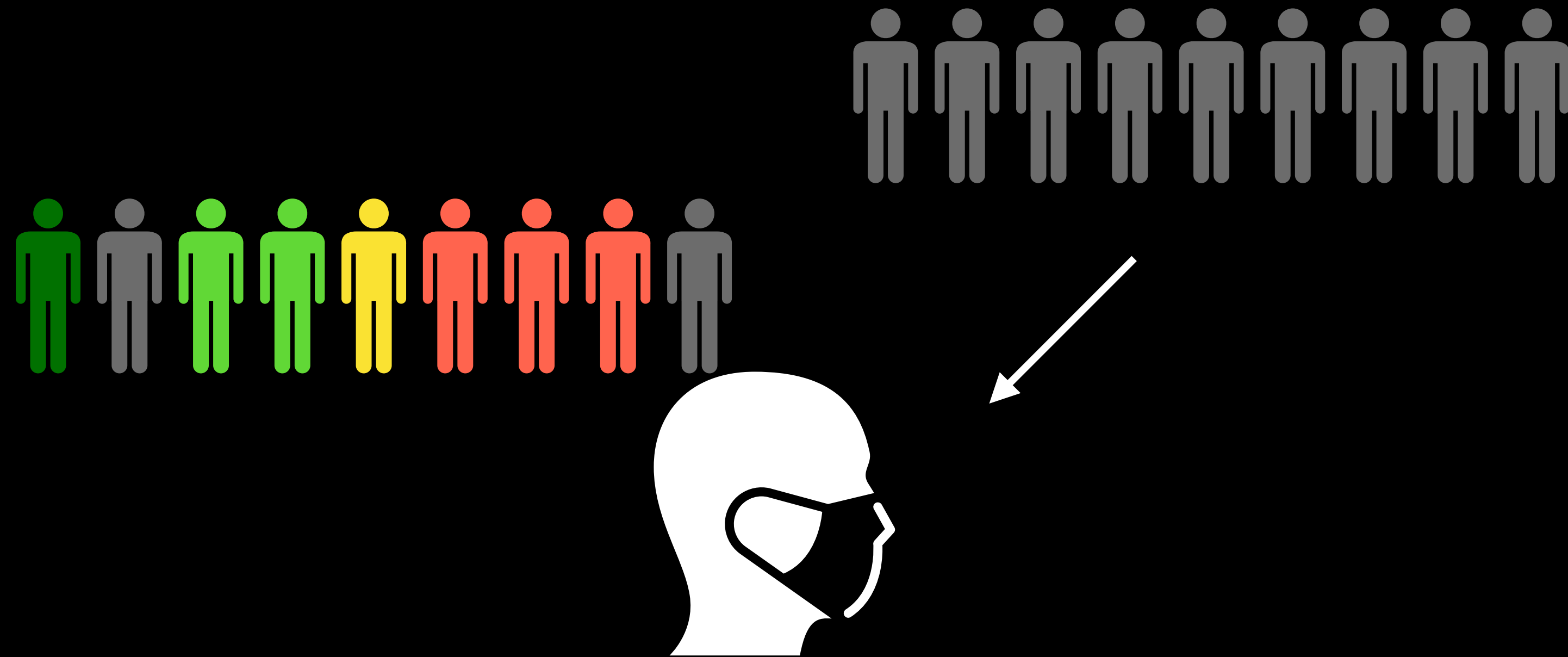
# Last but not least.



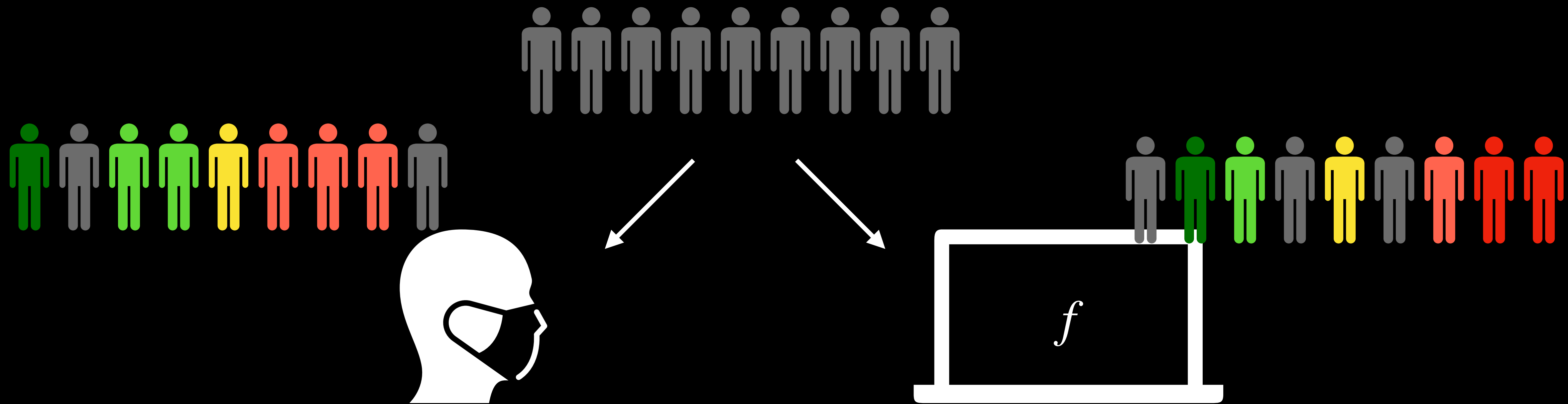
# Physicians and models function as a team in healthcare settings.



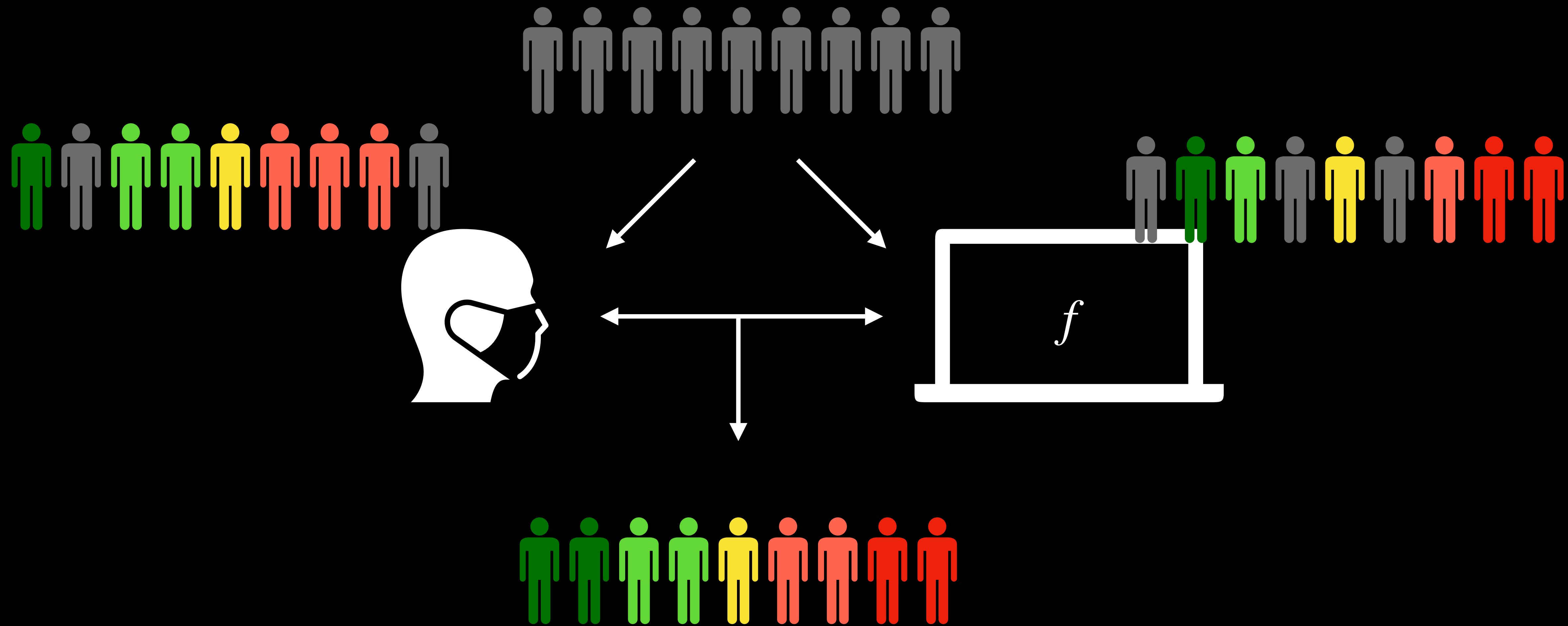
# Physicians and models function as a team in healthcare settings.



# Physicians and models function as a team in healthcare settings.



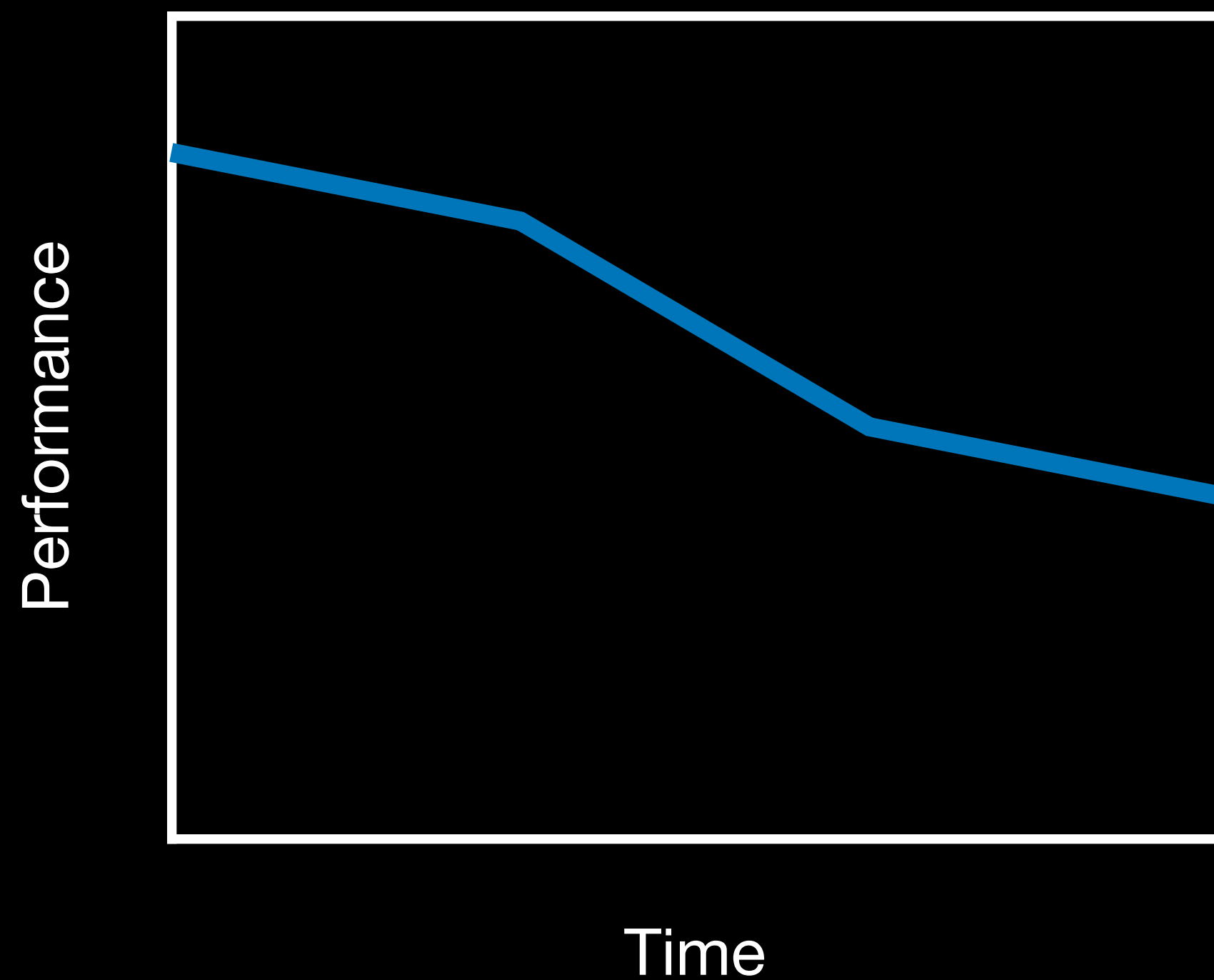
# Physicians and models function as a team in healthcare settings.



**This is complicated because we need to update models over time.**

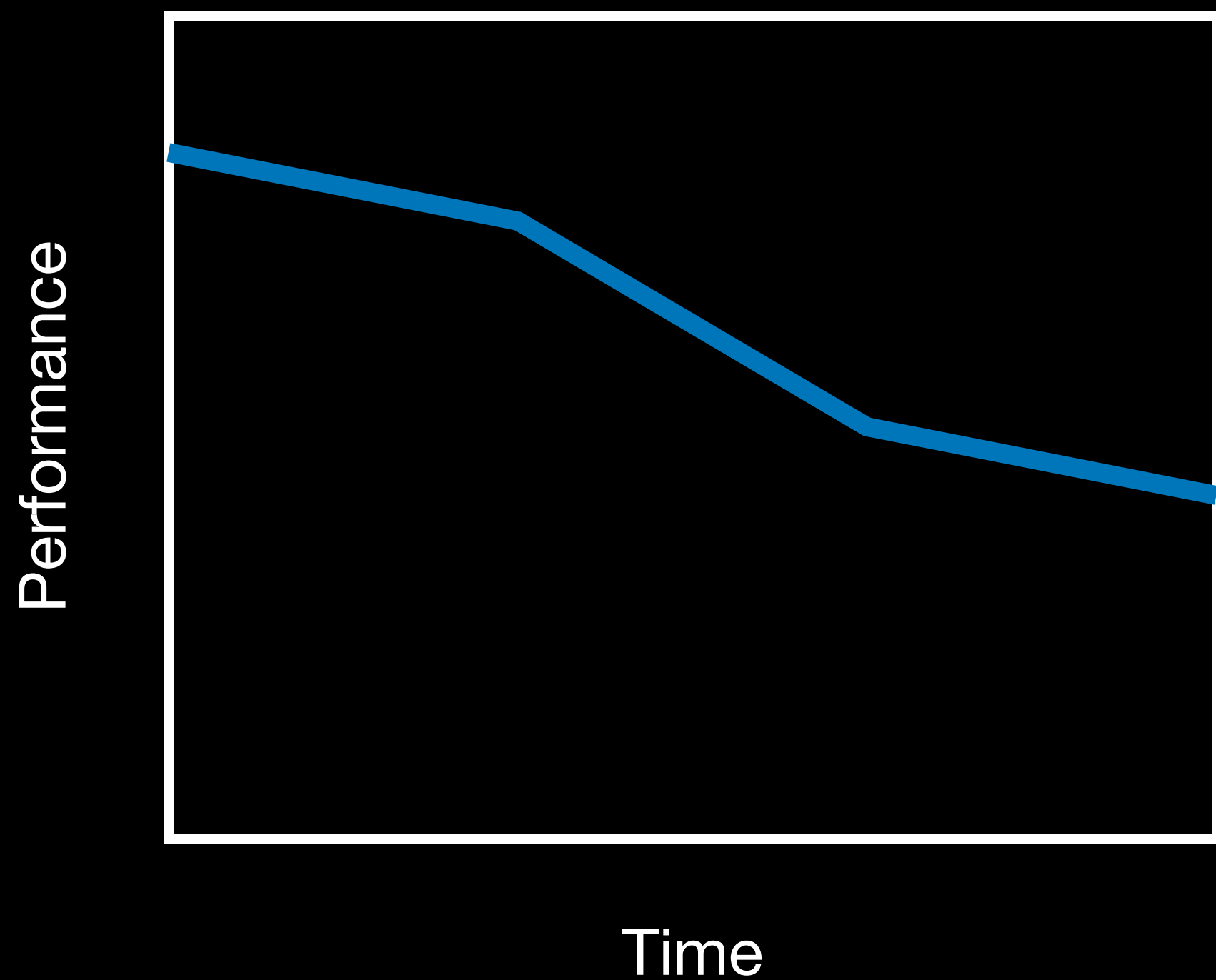
# This is complicated because we need to update models over time.

Counter performance degradation

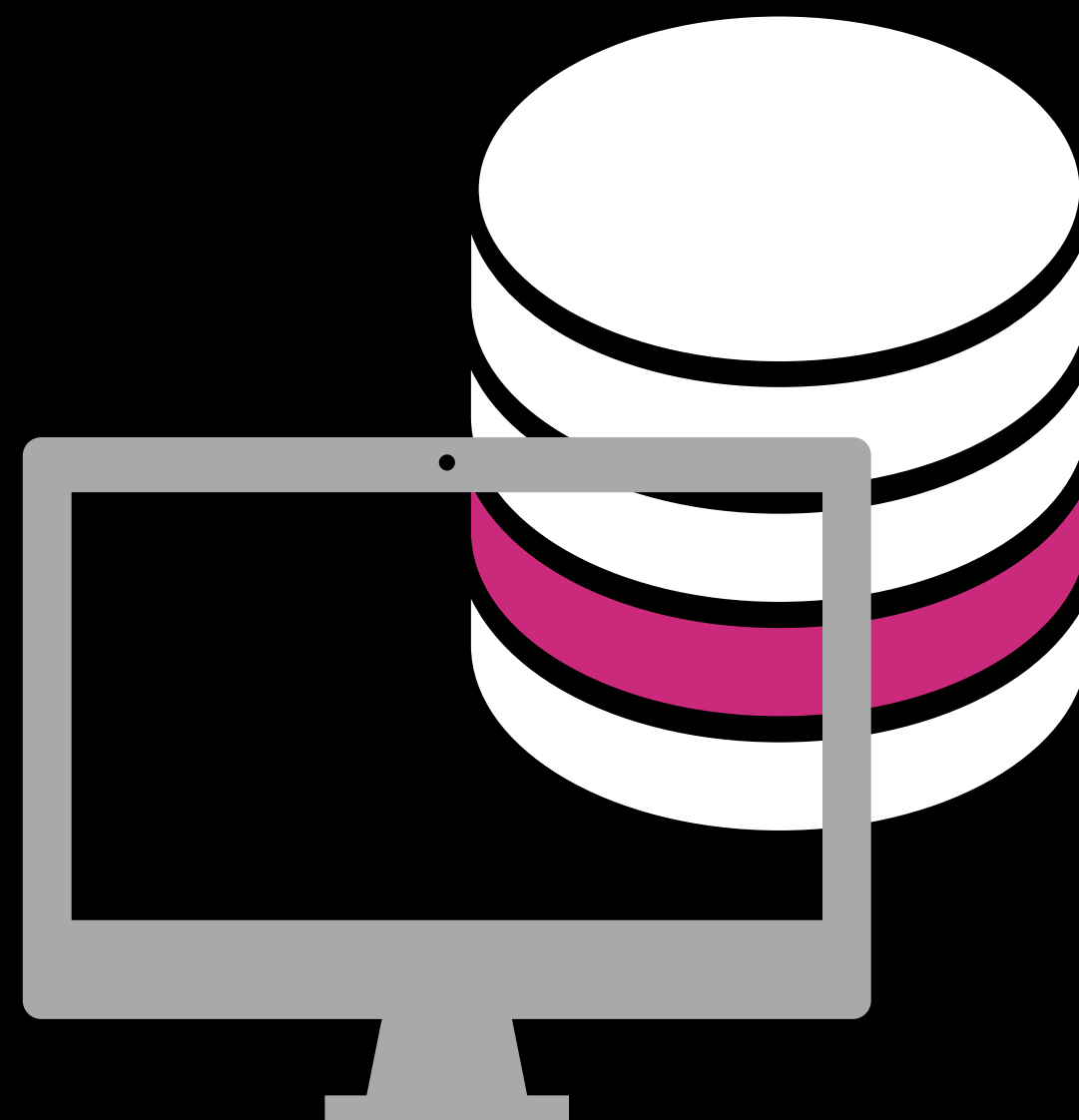


# This is complicated because we need to update models over time.

Counter performance degradation



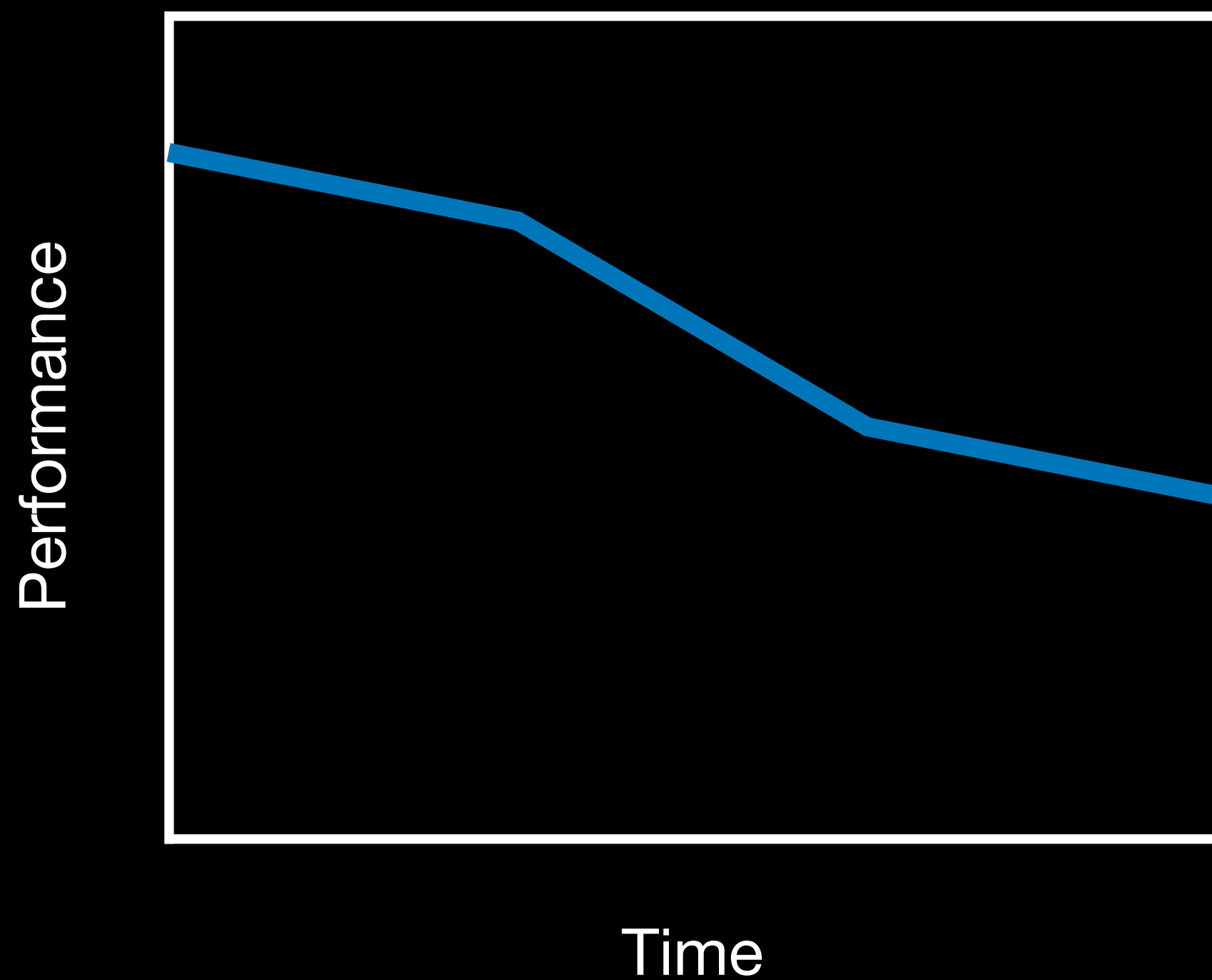
Adapt to infrastructure changes



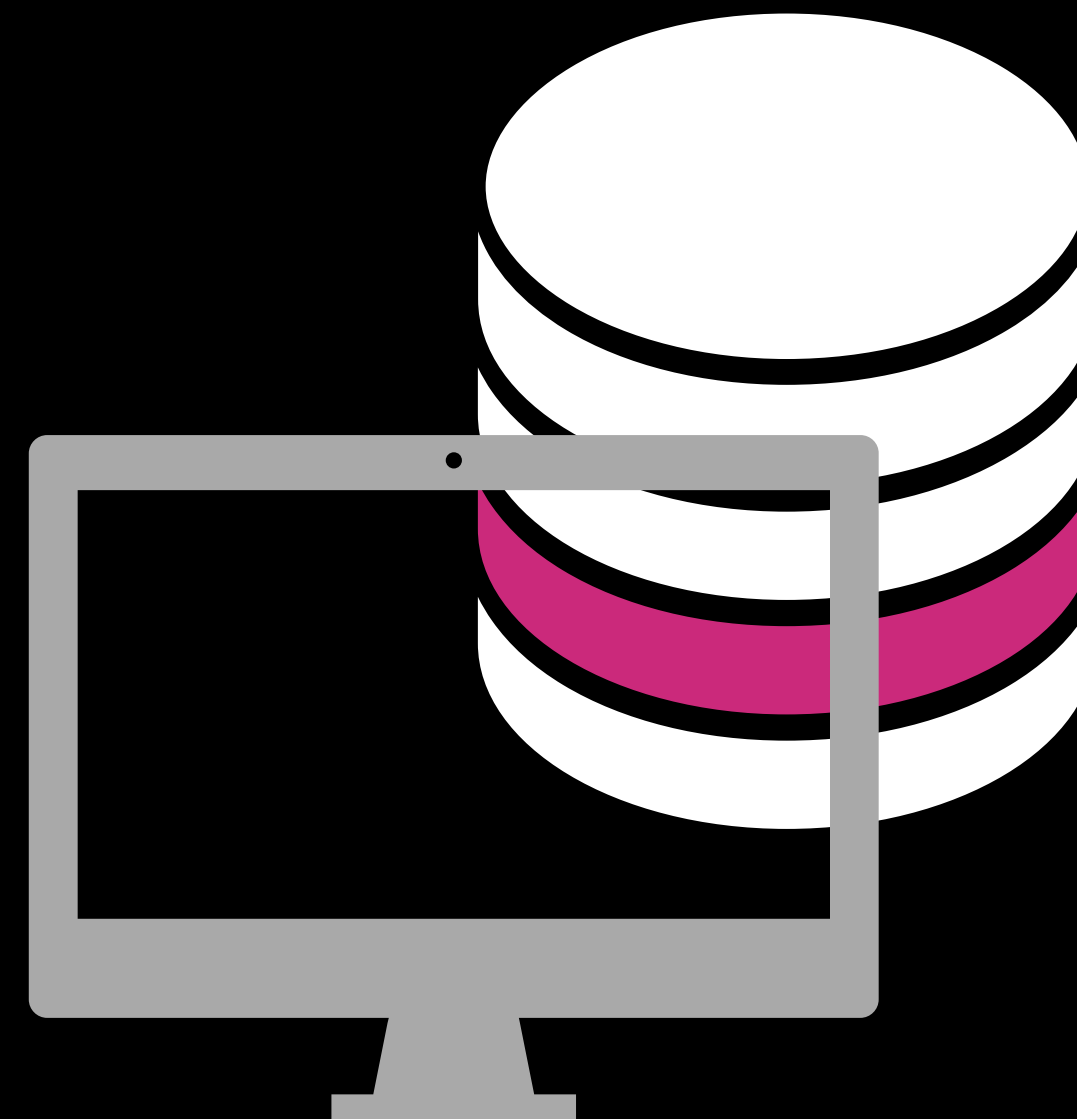


# This is complicated because we need to update models over time.

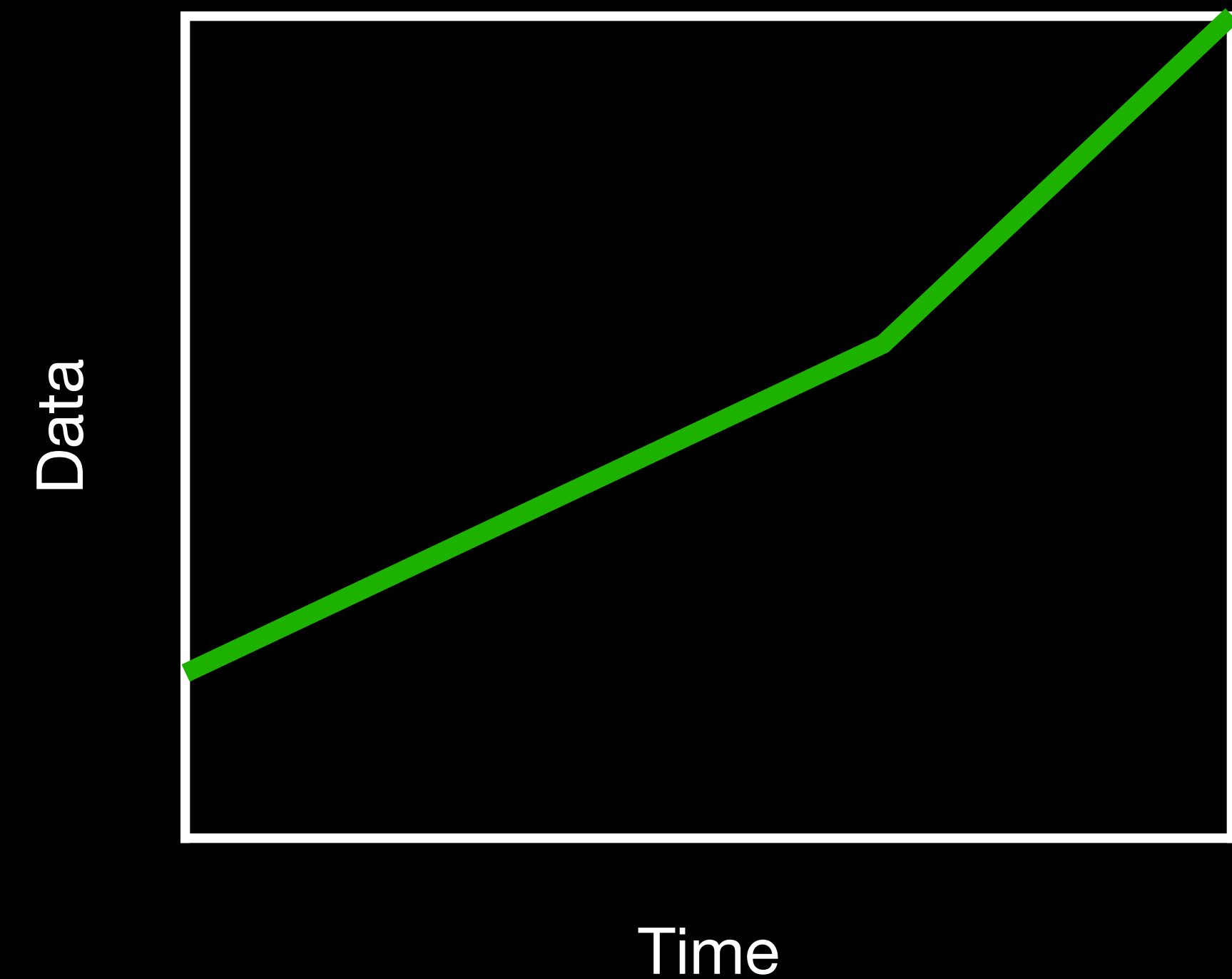
Counter performance degradation



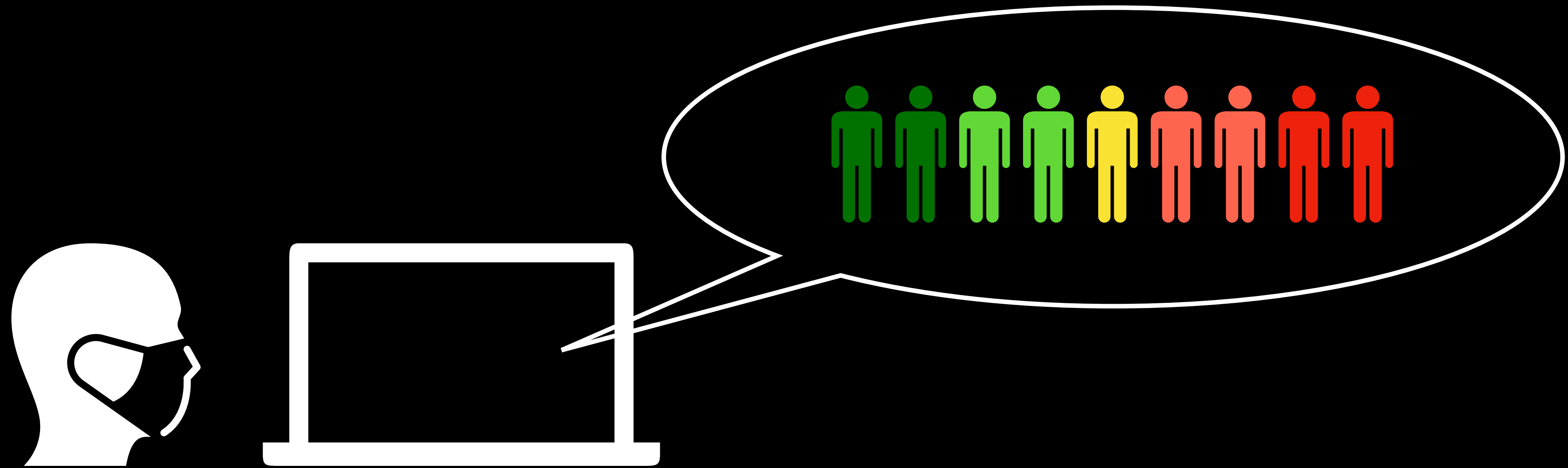
Adapt to infrastructure changes



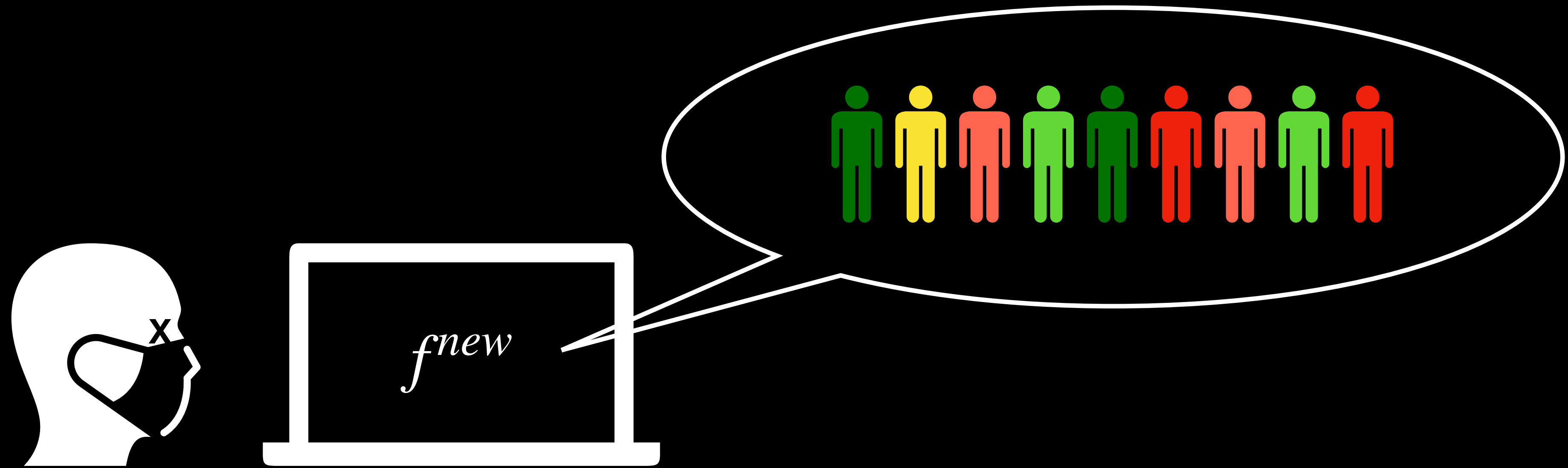
Incorporate new data



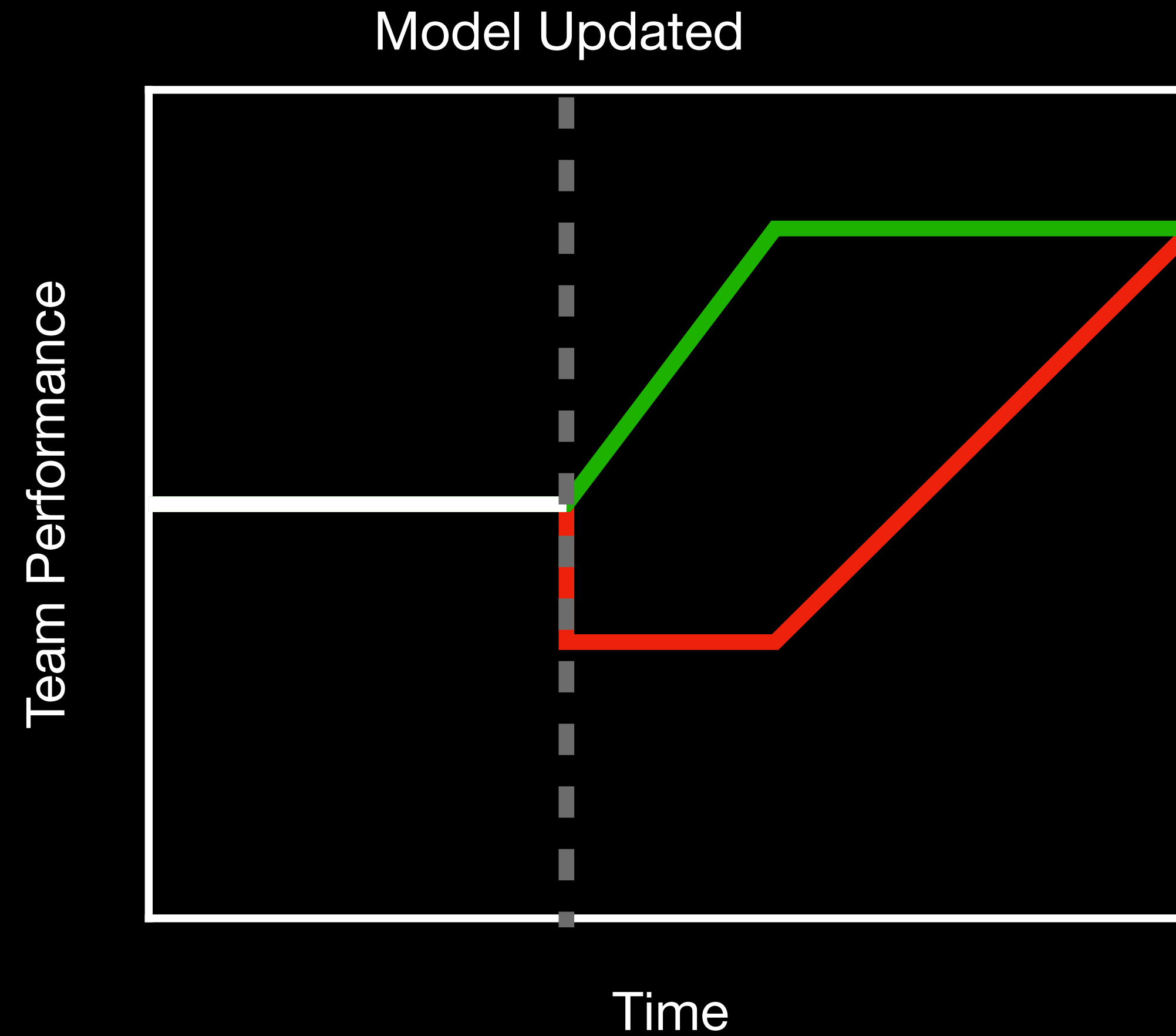
# Updates can mess with user expectations.



# Updates can mess with user expectations.



# Team performance may suffer if models don't meet user expectations.

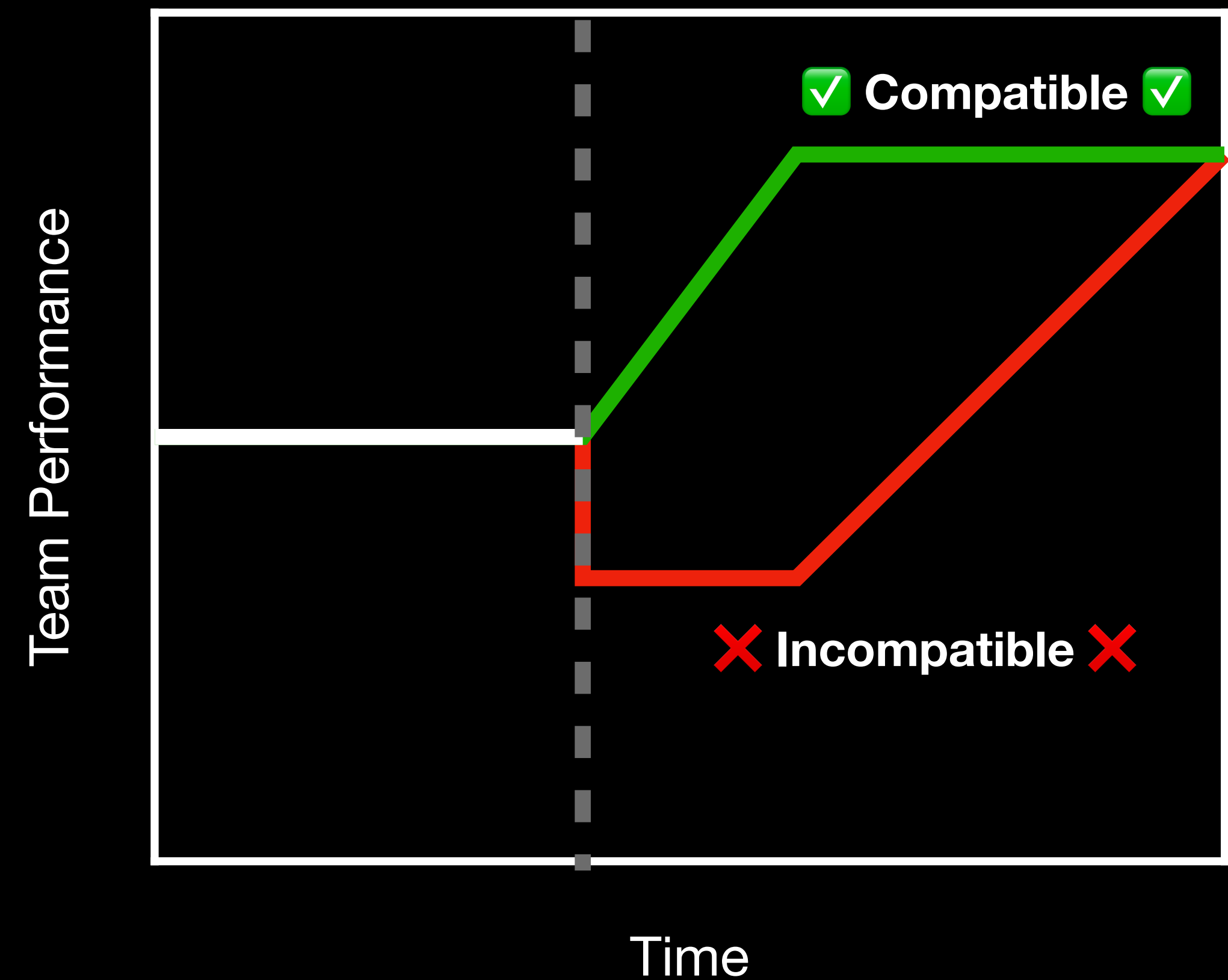


# Ideally updated models meet the expectations of users

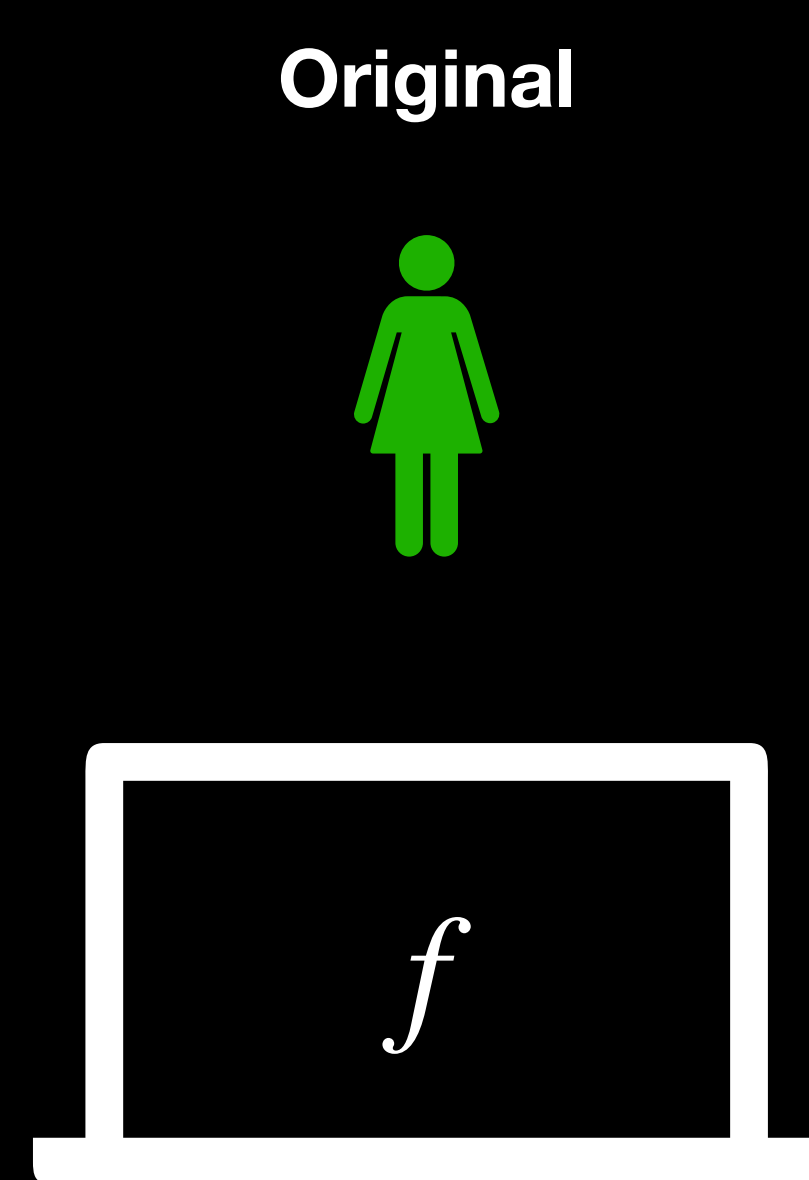
*Compatibility*: the amount an updated model continues the correct behavior of an original model

Way to measure user expectations

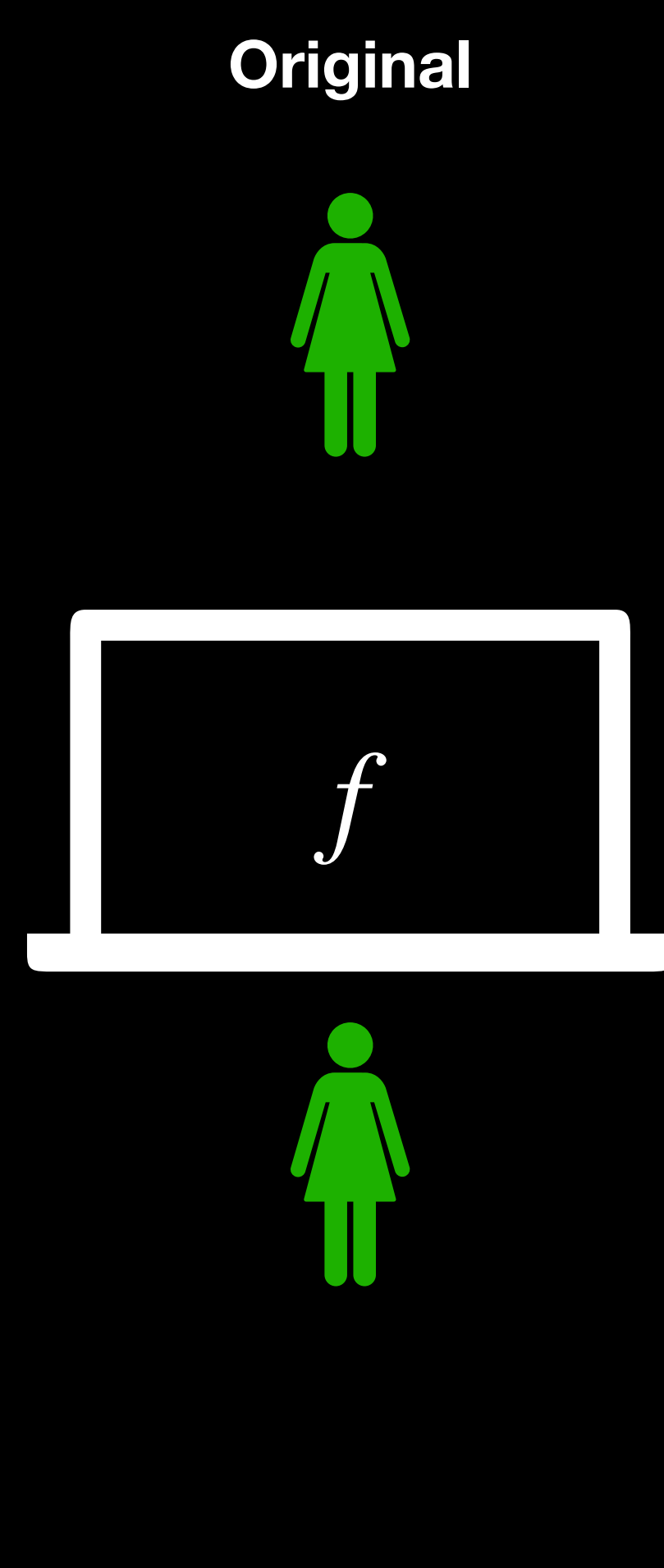
Goal: updated models should have high compatibility



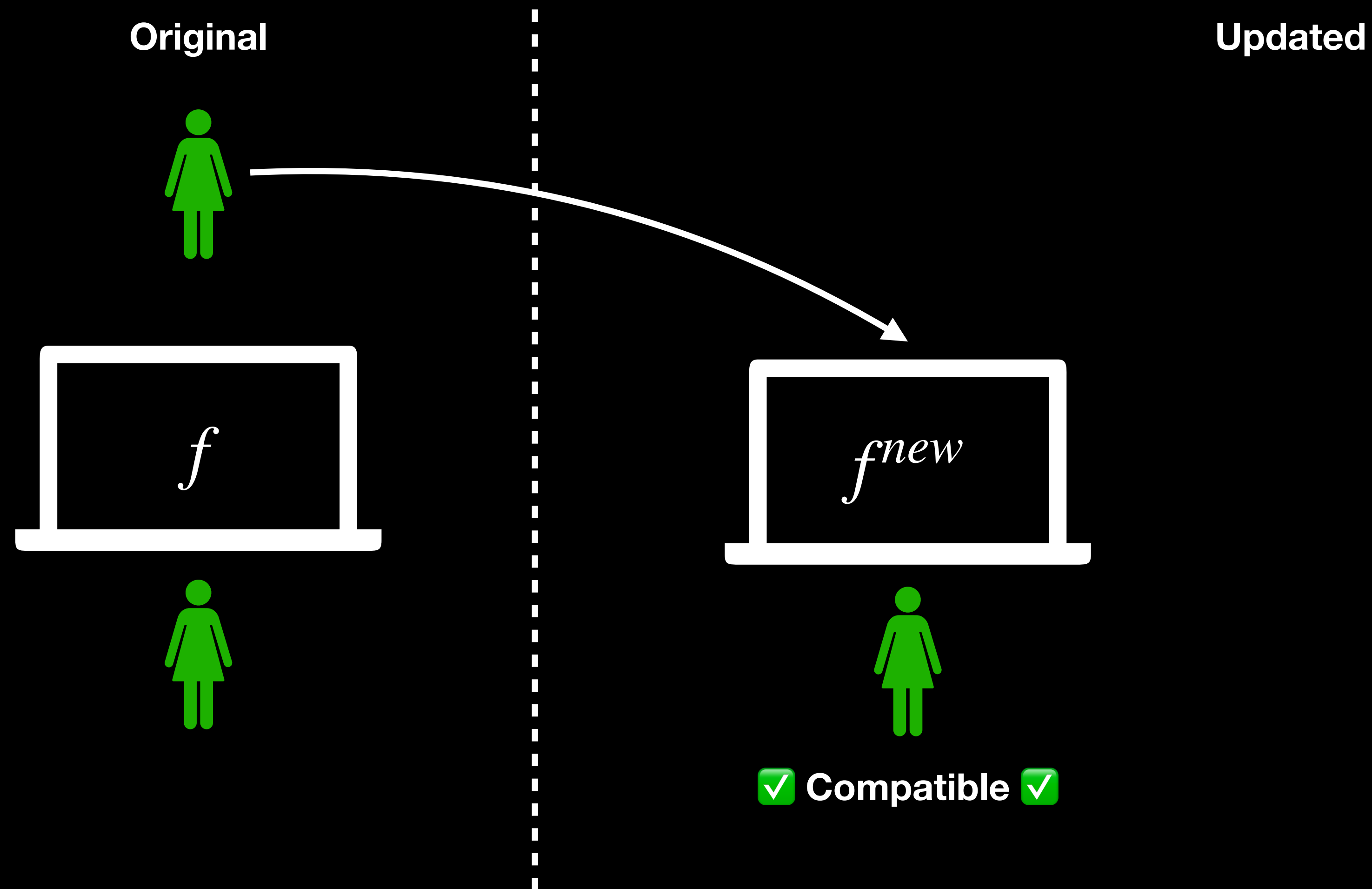
Compatibility can be assessed by using the original and updated models for the same predictive task.



Compatibility can be assessed by using the original and updated models for the same predictive task.

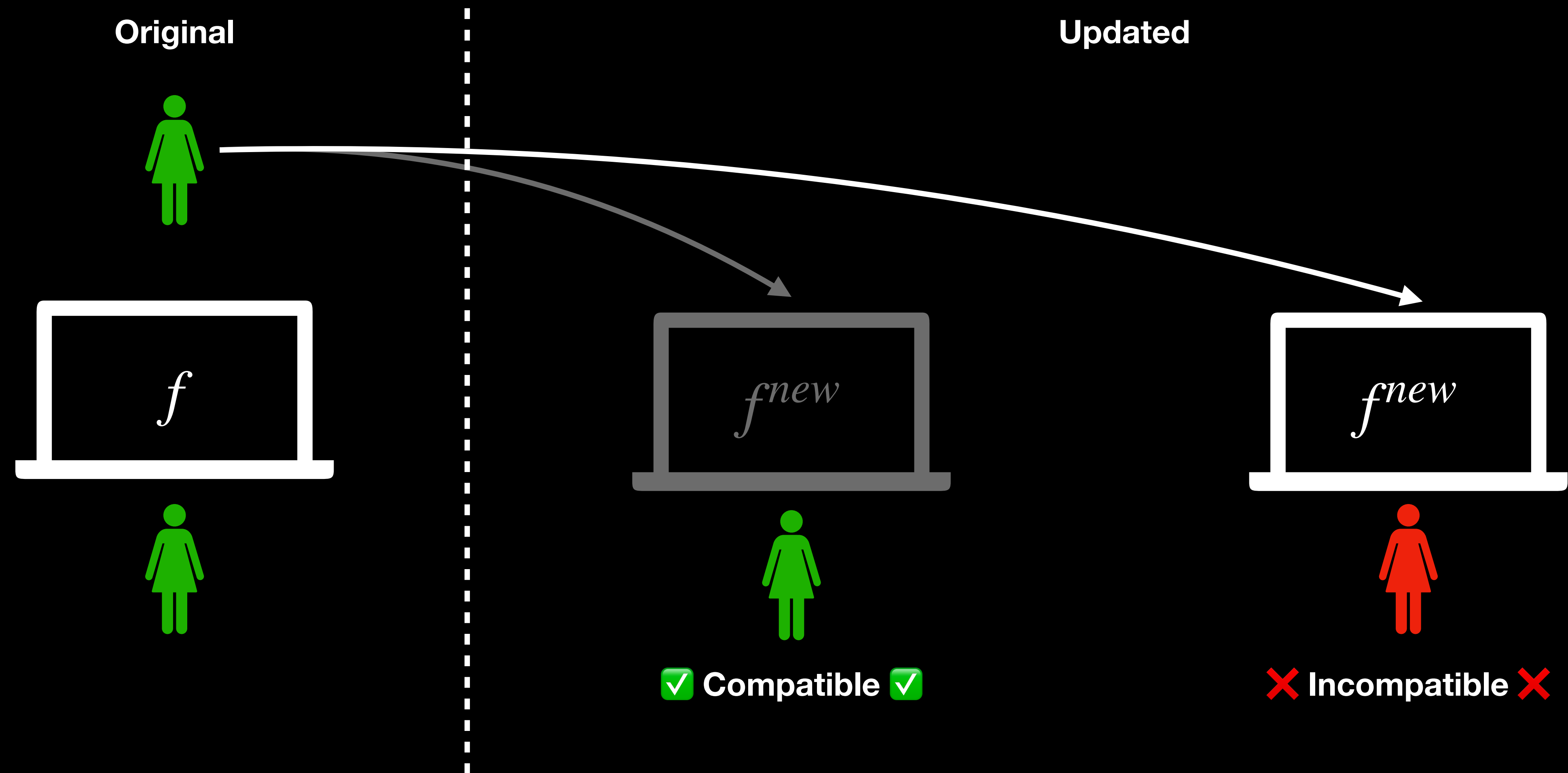


# Compatibility can be assessed by using the original and updated models for the same predictive task.





# Compatibility can be assessed by using the original and updated models for the same predictive task.



# Backwards Trust Compatibility $C^{BT}$

The chance that the updated model's labels are correct, given that the original model's labels were correct.

$$C^{BT}(f^o, f^u) = \frac{\# \text{ patients both models label correctly}}{\# \text{ patients original model labels correctly}}$$

# Backwards Trust Compatibility $C^{BT}$

The chance that the updated model's labels are correct, given that the original model's labels were correct.

$$C^{BT}(f^o, f^u) = \frac{\# \text{ patients both models label correctly}}{\# \text{ patients original model labels correctly}}$$

original model                  updated model

# Backwards Trust Compatibility $C^{BT}$

The chance that the updated model's labels are correct, given that the original model's labels were correct.

$$C^{BT}(f^o, f^u) = \frac{\text{\# patients both models label correctly}}{\text{\# patients original model labels correctly}}$$

# Backwards Trust Compatibility $C^{BT}$

The chance that the updated model's labels are correct, given that the original model's labels were correct.

$$C^{BT}(f^o, f^u) = \frac{\# \text{ patients both models label correctly}}{\# \text{ patients original model labels correctly}}$$

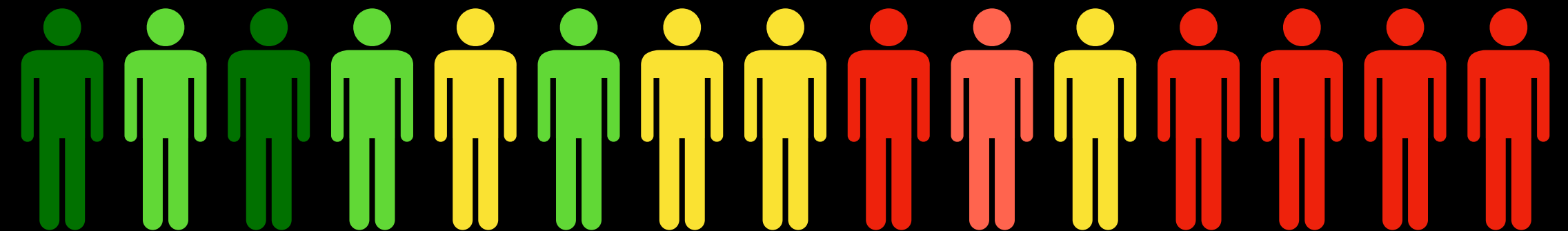
1  $\rightarrow$  perfect compatibility, 0  $\rightarrow$  perfect incompatibility

# Problems with existing compatibility measures.

Existing measure depends on equality comparison

Problematic for use in risk stratification model & healthcare settings

Depends on setting a single decision threshold

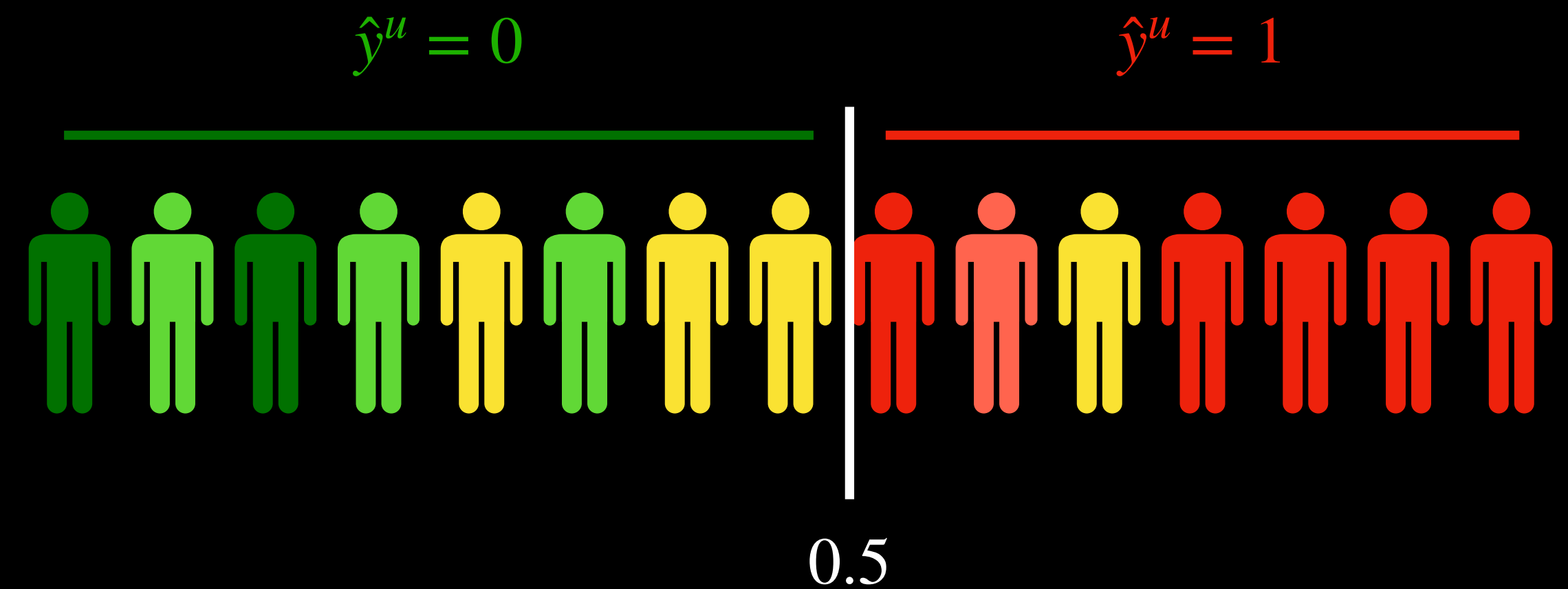


# Problems with existing compatibility measures.

Existing measure depends on equality comparison

Problematic for use in risk stratification model & healthcare settings

Depends on setting a single decision threshold



# Problems with existing compatibility measures.

Existing measure depends on equality comparison

Problematic for use in risk stratification model & healthcare settings

Depends on setting a single decision threshold



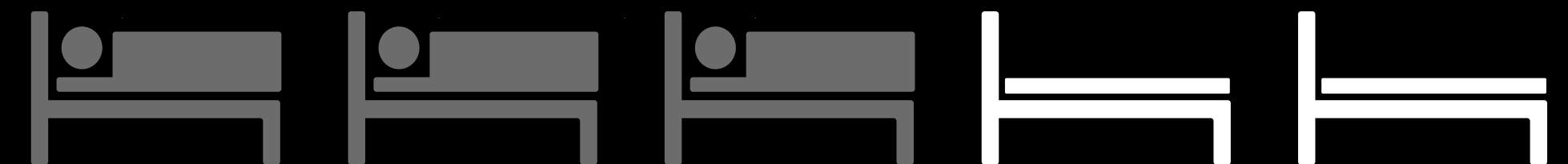


# Problems with existing compatibility measures.

Existing measure depends on equality comparison

Problematic for use in risk stratification model & healthcare settings

Depends on setting a single decision threshold



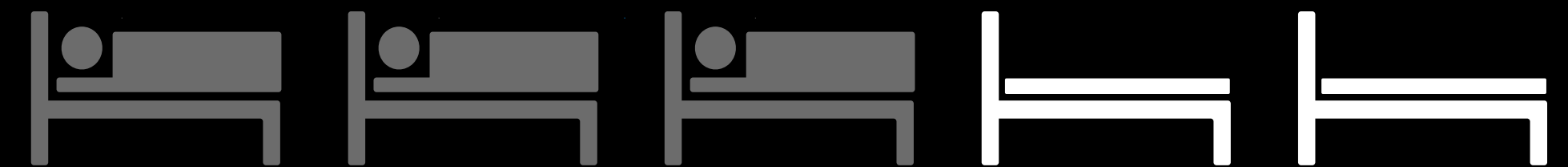
# Problems with existing compatibility measures.

Existing measure depends on equality comparison

Problematic for use in risk stratification model & healthcare settings

Depends on setting a single decision threshold

No direct relationship with AUROC



# Our contributions

Define a new rank-based compatibility measure ( $C^R$ )

Characterize  $C^R$  and its relationship with AUROC

Custom loss function to engineer model updates with improved  $C^R$

Intuition:  $C^R$  should inherit from both  $C^{BT}$  & AUROC

$$C^{BT}(f^o, f^u) = \frac{\text{\# patients both models label correctly}}{\text{\# patients original model labels correctly}}$$

$$AUROC(f^o) = \frac{\sum_{i \in I^0} \sum_{j \in I^1} \mathbf{1}(\hat{p}_i^o < \hat{p}_j^o)}{m}$$

Evaluate correct behavior of both models  
Normalized based on original model's behavior

“Correctness” based on risk estimate ordering

# Rank-based compatibility $C^R$

Agreement of risk estimate rankings produced by original & updated models given original ranked correctly:

$$C^R(f^o, f^u) = \frac{\sum_{i \in I^0} \sum_{j \in I^1} \mathbf{1}(\hat{p}_i^o < \hat{p}_j^o) \cdot \mathbf{1}(\hat{p}_i^u < \hat{p}_j^u)}{\sum_{i \in I^0} \sum_{j \in I^1} \mathbf{1}(\hat{p}_i^o < \hat{p}_j^o)}$$

# Rank-based compatibility $C^R$

Agreement of risk estimate rankings produced by original & updated models given original ranked correctly:

$$C^R(f^o, f^u) = \frac{\sum_{i \in I^0} \sum_{j \in I^1} \mathbf{1}(\hat{p}_i^o < \hat{p}_j^o) \cdot \mathbf{1}(\hat{p}_i^u < \hat{p}_j^u)}{\sum_{i \in I^0} \sum_{j \in I^1} \mathbf{1}(\hat{p}_i^o < \hat{p}_j^o)}$$

original model      updated model

# Rank-based compatibility $C^R$

Agreement of risk estimate rankings produced by original & updated models given original ranked correctly:

$$C^R(f^o, f^u) = \frac{\sum_{i \in I^0} \sum_{j \in I^1} \mathbf{1}(\hat{p}_i^o < \hat{p}_j^o) \cdot \mathbf{1}(\hat{p}_i^u < \hat{p}_j^u)}{\sum_{i \in I^0} \sum_{j \in I^1} \mathbf{1}(\hat{p}_i^o < \hat{p}_j^o)}$$

original model  
ranks correctly

# Rank-based compatibility $C^R$

Agreement of risk estimate rankings produced by original & updated models given original ranked correctly:

$$C^R(f^o, f^u) = \frac{\sum_{i \in I^0} \sum_{j \in I^1} \mathbf{1}(\hat{p}_i^o < \hat{p}_j^o) \cdot \mathbf{1}(\hat{p}_i^u < \hat{p}_j^u)}{\sum_{i \in I^0} \sum_{j \in I^1} \mathbf{1}(\hat{p}_i^o < \hat{p}_j^o)}$$

updated model ranks correctly

original model ranks correctly



# Rank-based compatibility $C^R$

Agreement of risk estimate rankings produced by original & updated models given original ranked correctly:

$$C^R(f^o, f^u) = \frac{\sum_{i \in I^0} \sum_{j \in I^1} \mathbf{1}(\hat{p}_i^o < \hat{p}_j^o) \cdot \mathbf{1}(\hat{p}_i^u < \hat{p}_j^u)}{\sum_{i \in I^0} \sum_{j \in I^1} \mathbf{1}(\hat{p}_i^o < \hat{p}_j^o)}$$

both models rank correctly

original model ranks correctly

# Rank-based compatibility $C^R$

Agreement of risk estimate rankings produced by original & updated models given original ranked correctly:

$$C^R(f^o, f^u) = \frac{\sum_{i \in I^0} \sum_{j \in I^1} \mathbf{1}(\hat{p}_i^o < \hat{p}_j^o) \cdot \mathbf{1}(\hat{p}_i^u < \hat{p}_j^u)}{\sum_{i \in I^0} \sum_{j \in I^1} \mathbf{1}(\hat{p}_i^o < \hat{p}_j^o)}$$

1 → perfect compatibility, 0 → perfect incompatibility

$C^R$  is a new compatibility measure inspired by AUROC

Not threshold dependent.

Has a direct relationship with AUROC which we can use to assess trade-offs.

Lower bound of  $C^R$  expected to increase as model performance increases.

**Q1: Do we get  $C^R$  for free when we make updated models targeting AUROC?**

More specifically: do we observe  $C^R = 1$  (or very close) when we train updated models using binary cross entropy loss?

Hypothesis: No, analytically we'd expect that  $C^R$  is centered at a region away from the upper and lower bounds.

# Q1: Experimental Setup

Original Model

$n = 1,000$

Original Model Development

$n = 500$

Original Model Validation

$n = 500$

Updated Model Dataset

$n = 5,000$

Updated Model Development

$n = 2,500$

Updated Model Validation

$n = 2,500$

Evaluation  
Dataset

$n = 2,577$

**MIMIC-III Mortality Dataset**

# Q1: Experimental Setup

Original Model  
 $n = 1,000$

Original Model Development  
 $n = 500$

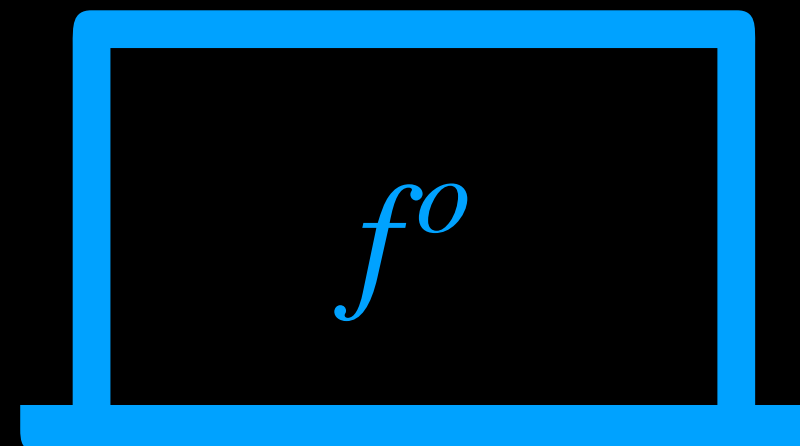
Original Model Validation  
 $n = 500$

Updated Model Dataset  
 $n = 5,000$

Updated Model Development  
 $n = 2,500$

Updated Model Validation  
 $n = 2,500$

Evaluation  
Dataset  
 $n = 2,577$

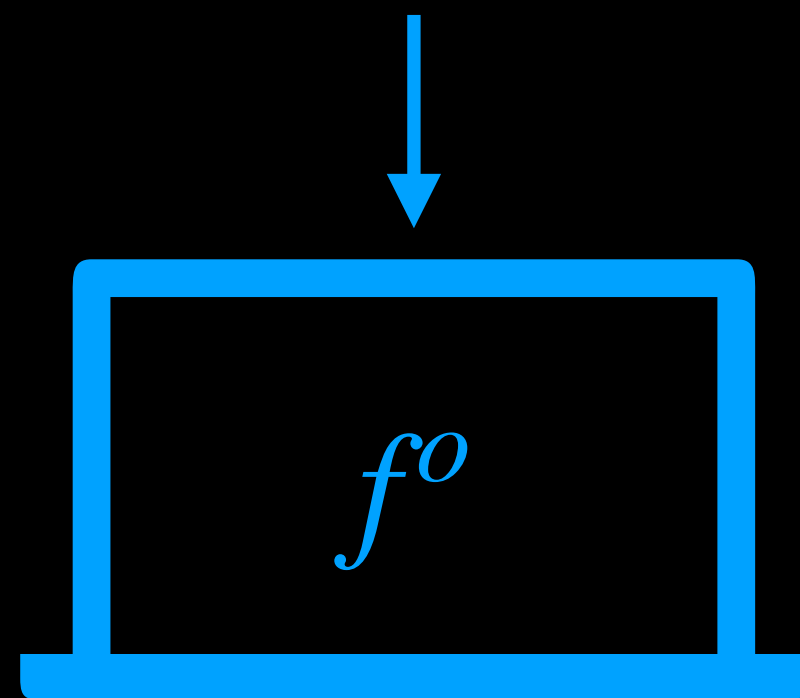


# Q1: Experimental Setup

Original Model  
 $n = 1,000$

Original Model Development  
 $n = 500$

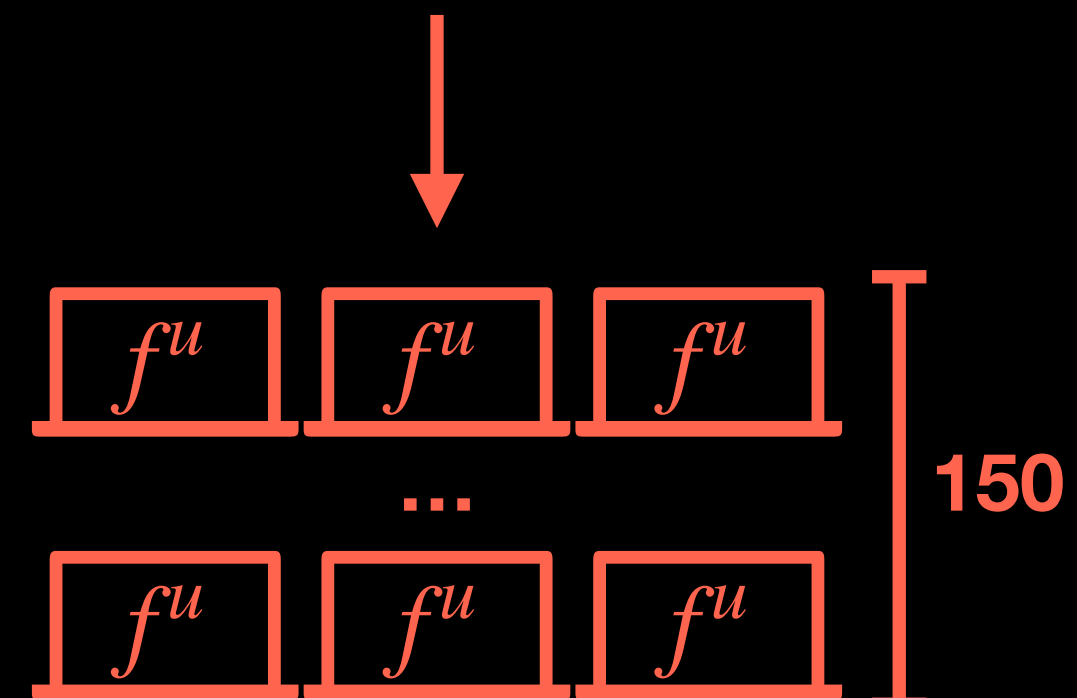
Original Model Validation  
 $n = 500$



Updated Model Dataset  
 $n = 5,000$

Updated Model Development  
 $n = 2,500$

Updated Model Validation  
 $n = 2,500$



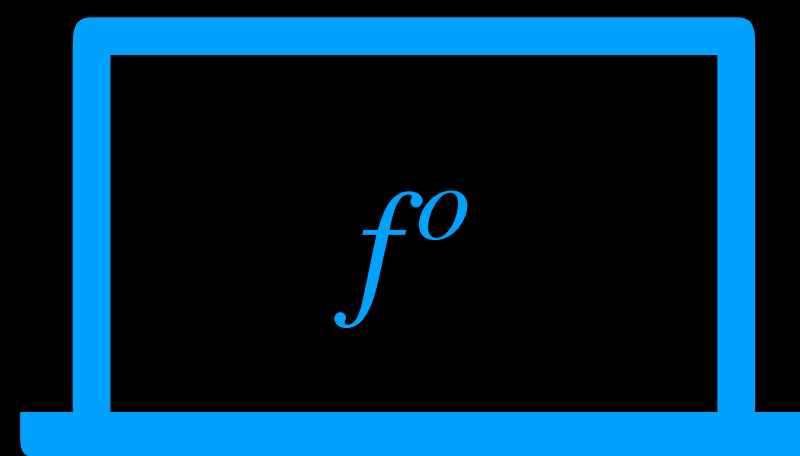
Evaluation  
Dataset  
 $n = 2,577$

# Q1: Experimental Setup

Original Model  
 $n = 1,000$

Original Model Development  
 $n = 500$

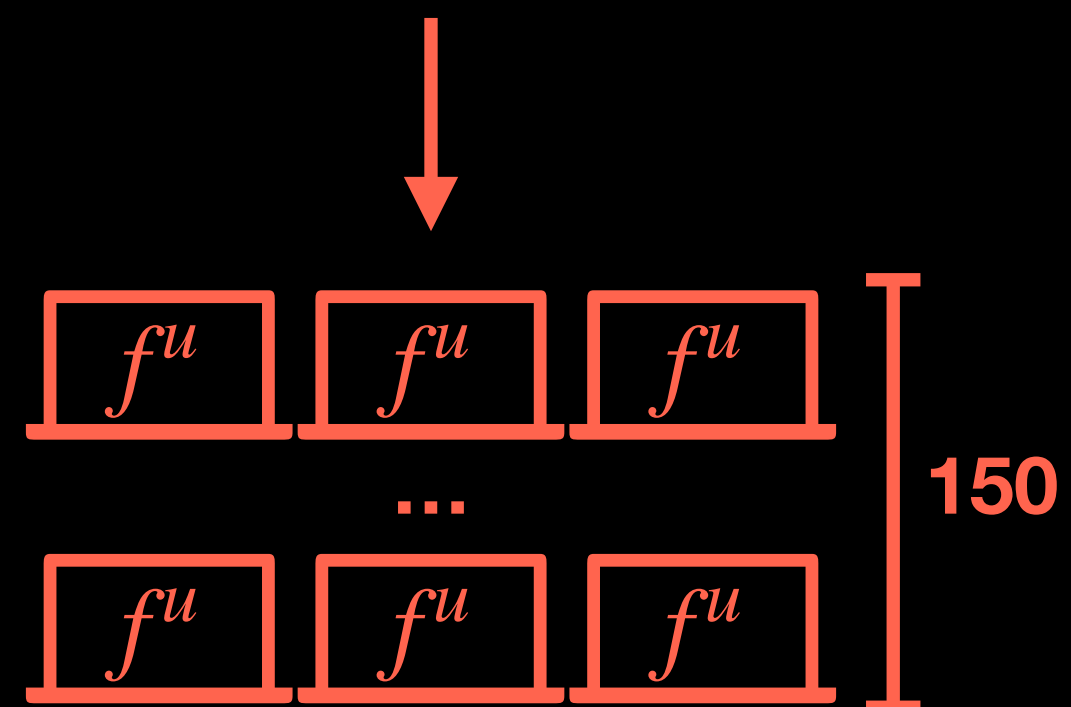
Original Model Validation  
 $n = 500$



Updated Model Dataset  
 $n = 5,000$

Updated Model Development  
 $n = 2,500$

Updated Model Validation  
 $n = 2,500$

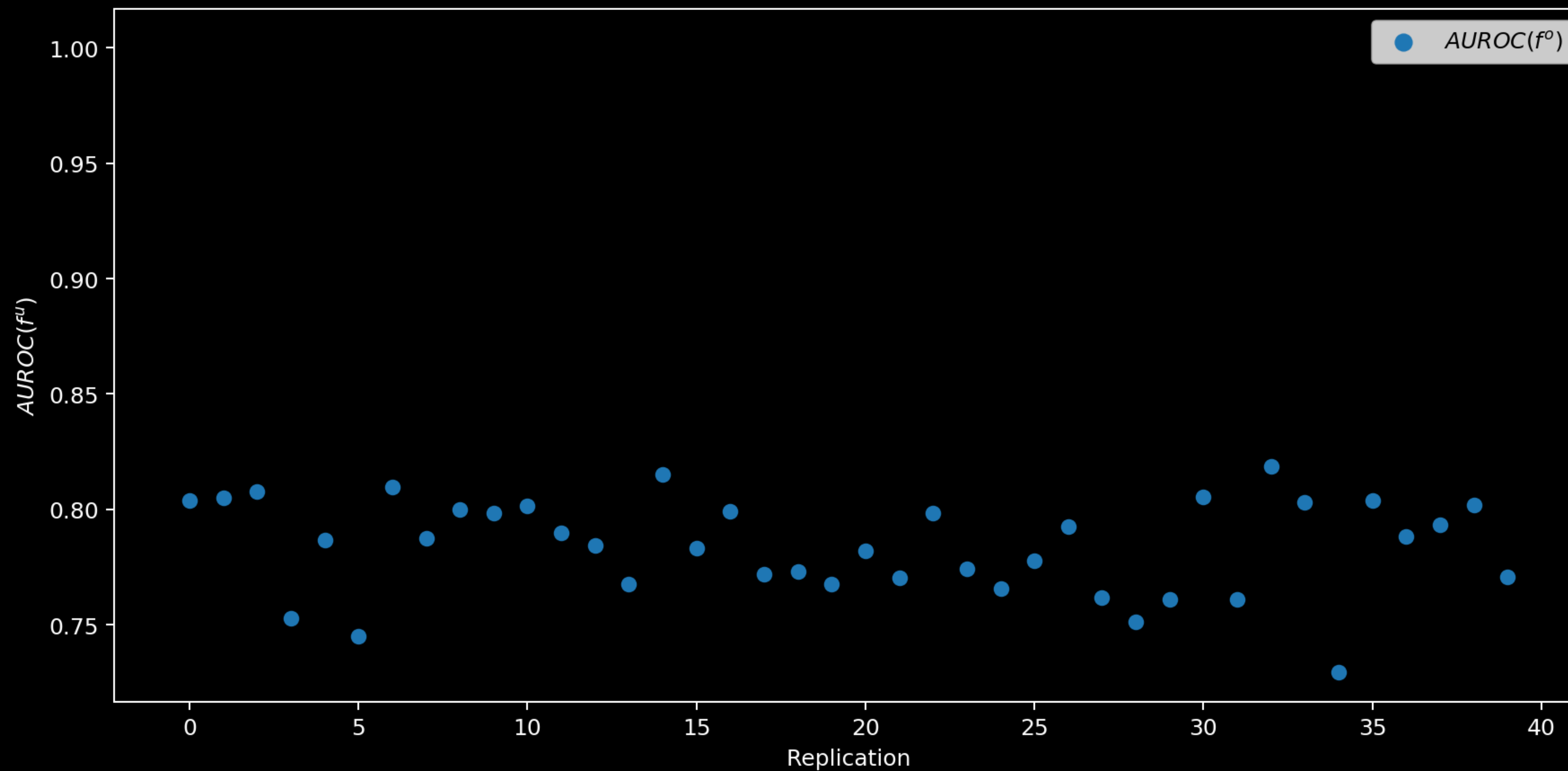


Evaluation  
Dataset  
 $n = 2,577$

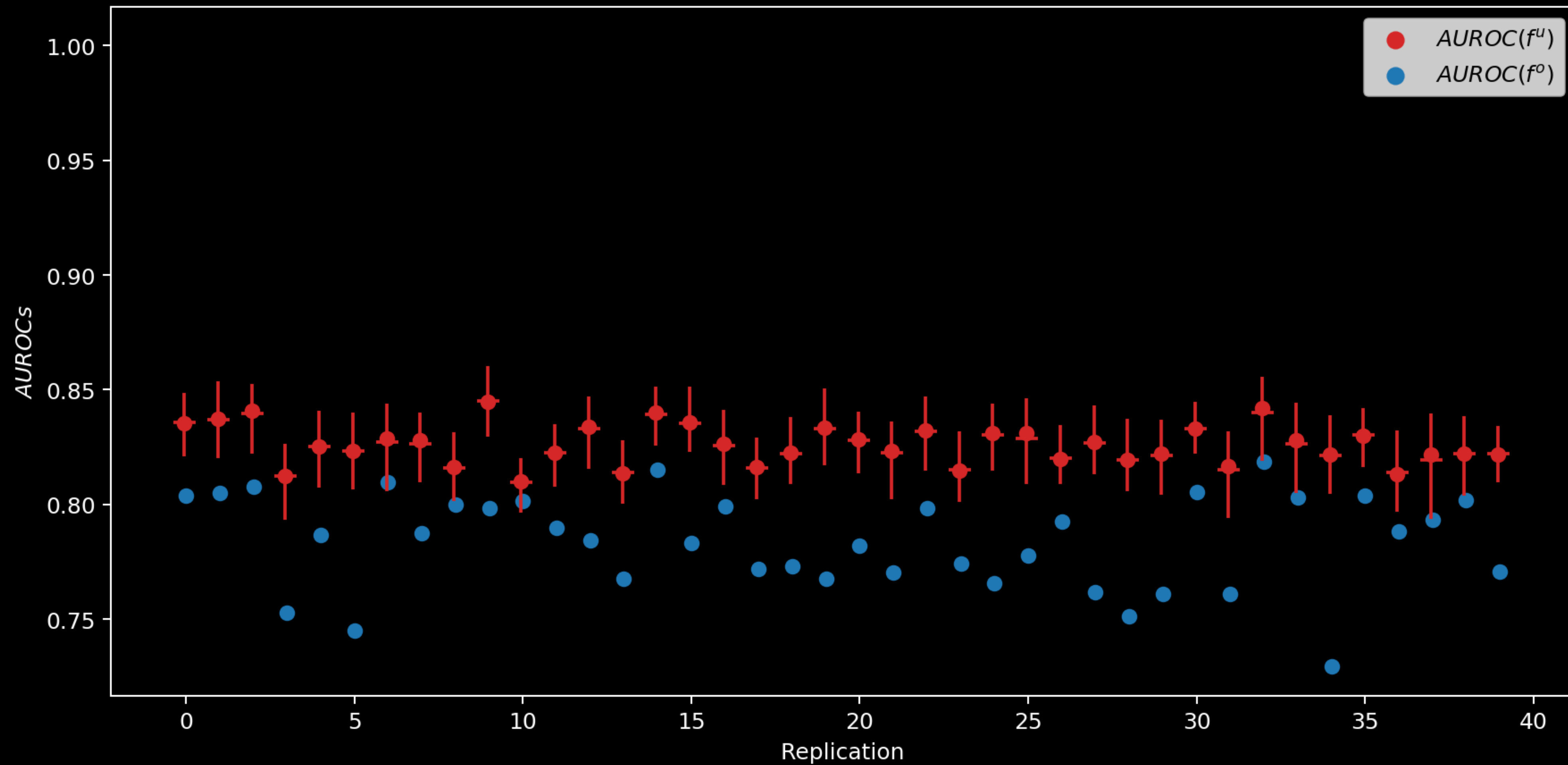
$AUROC(f^o)$   
 $AUROC(f^u)$   
 $C^R(f^o, f^u)$



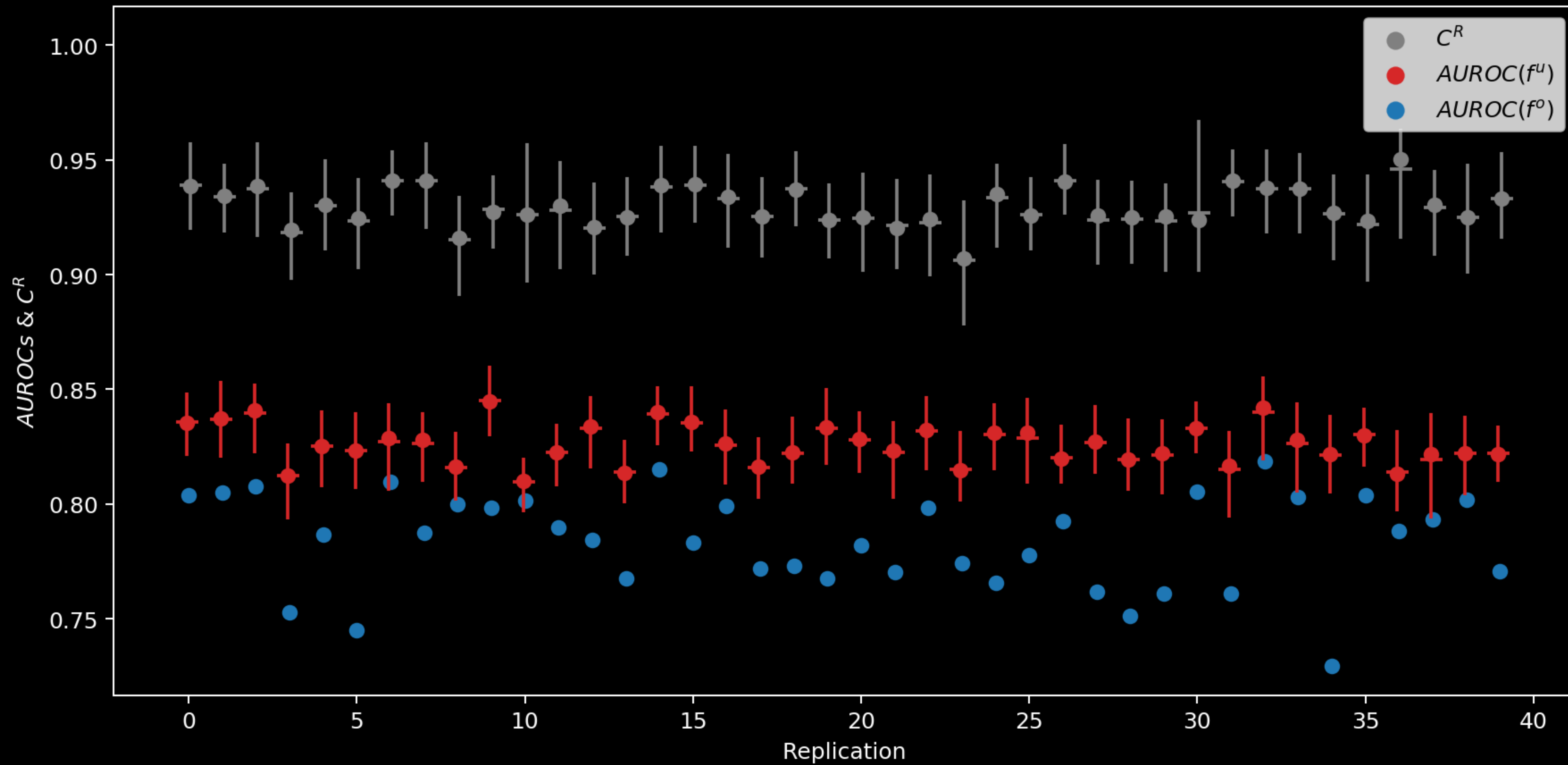
# Results



# Results



# Results



**Q1: Do we get  $C^R$  for free when we make updated models targeting AUROC?**

No.

We observe updated models have a limited range in  $C^R$ .

Analytical results suggest that there's a large search space.

Motivates techniques to search for updated models that have higher  $C^R$ .

# Risk stratification models are usually trained with binary cross entropy loss.

Binary cross entropy loss:

$$\mathcal{L}^{BCE}(f) = - \sum_{i \in I^0} \log(1 - \hat{p}_i) - \sum_{j \in I^1} \log(\hat{p}_j)$$

Minimization of  $\mathcal{L}^{BCE}$  leads to higher AUROC because risk estimates tend to align with labels.

# Risk stratification models are usually trained with binary cross entropy loss.

Binary cross entropy loss:

$$\mathcal{L}^{BCE}(f) = - \sum_{i \in I^0} \log(1 - \hat{p}_i) - \sum_{j \in I^1} \log(\hat{p}_j)$$

Make 0-labeled patients have low risk estimates

Minimization of  $\mathcal{L}^{BCE}$  leads to higher AUROC because risk estimates tend to align with labels.

# Risk stratification models are usually trained with binary cross entropy loss.

Binary cross entropy loss:

$$\mathcal{L}^{BCE}(f) = - \sum_{i \in I^0} \log(1 - \hat{p}_i) - \sum_{j \in I^1} \log(\hat{p}_j)$$

Make 0-labeled patients have low risk estimates

Make 1-labeled patients have high risk estimates

Minimization of  $\mathcal{L}^{BCE}$  leads to higher AUROC because risk estimates tend to align with labels.

# Risk stratification models are usually trained with binary cross entropy loss.

Binary cross entropy loss:

$$\mathcal{L}^{BCE}(f) = - \sum_{i \in I^0} \log(1 - \hat{p}_i) - \sum_{j \in I^1} \log(\hat{p}_j)$$

Make 0-labeled patients have low risk estimates

Make 1-labeled patients have high risk estimates

Minimization of  $\mathcal{L}^{BCE}$  leads to higher AUROC because risk estimates tend to align with labels.

No focus on compatibility between the updated and original model.



# We introduce rank-based incompatibility loss.

Rank-based incompatibility loss:

$$\mathcal{L}^R(f^o, f^u) = 1 - C^R(f^o, f^u)$$

Minimization of  $\mathcal{L}^R$  will lead to higher levels of  $C^R$ .

Differentiable approximation  $\widetilde{\mathcal{L}}^R$  for SGD.

# Weighted loss trades-off between binary cross entropy and compatibility.

Weighted loss function:

$$\alpha \mathcal{L}^{BCE}(f^u) + (1 - \alpha) \widetilde{\mathcal{L}}^R(f^o, f^u)$$

where  $\alpha \in [0,1]$

When:

$\alpha = 1$  then only minimize  $\mathcal{L}^{BCE}$ ,  $\uparrow$  AUROC

$\alpha = 0$  then only minimize  $\widetilde{\mathcal{L}}^R$ ,  $\uparrow C^R$

$\alpha = 0.5$  then balance  $\mathcal{L}^{BCE}$  and  $\widetilde{\mathcal{L}}^R$

**Q2: Can we make updated models with higher levels of  $C^R$ ?**

Specifically: Compared to standard update model generation and selection approaches, can we use the weighted loss function to generate updates with better  $C^R$ ?

Hypothesis: using weighted loss function will produce models with better  $C^R$ .

*Also, can this be accomplished without a loss of AUROC?*

# Q2: Extends Previous Experimental Setup

Original Model  
 $n = 1,000$

Original Model Development  
 $n = 500$

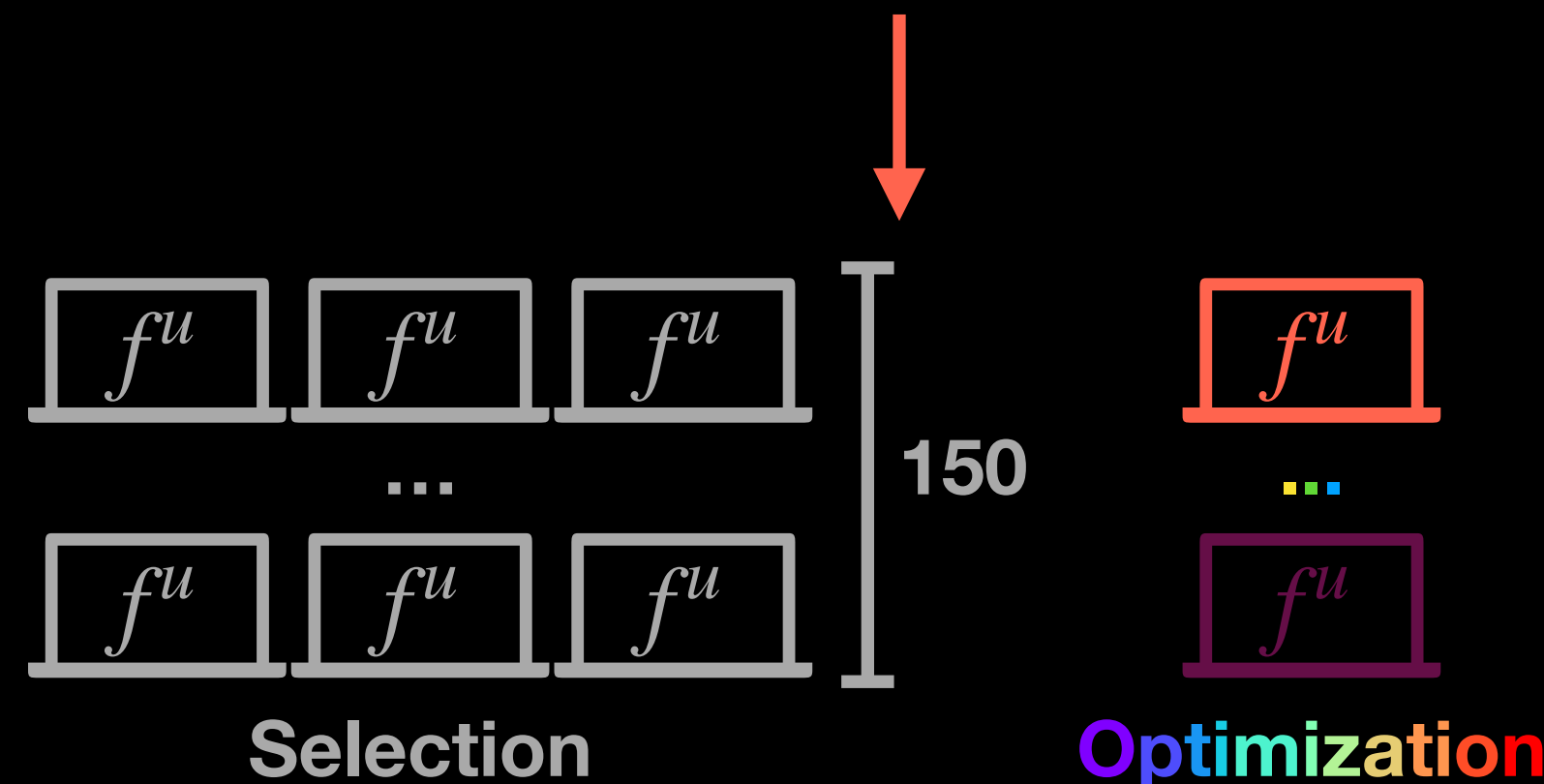
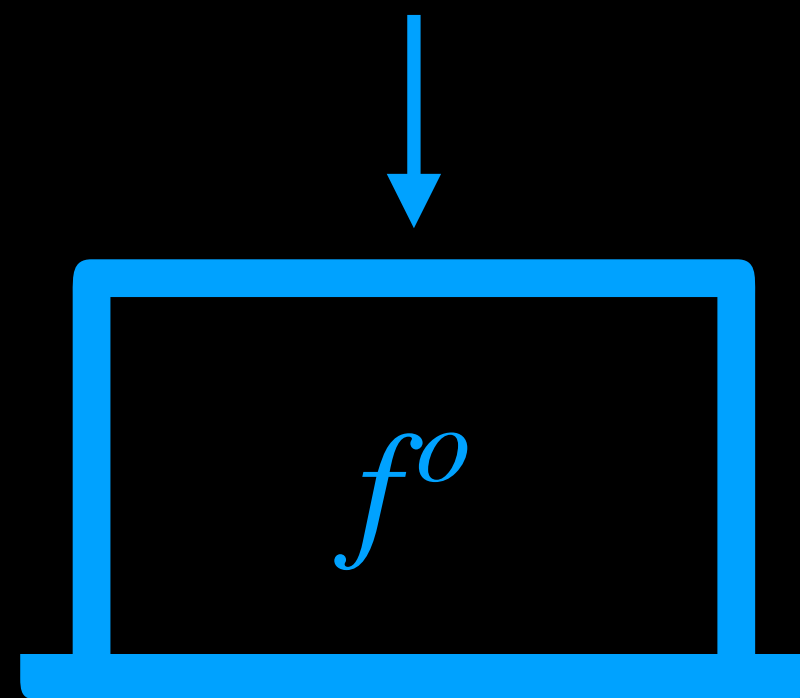
Original Model Validation  
 $n = 500$

Updated Model Dataset  
 $n = 5,000$

Updated Model Development  
 $n = 2,500$

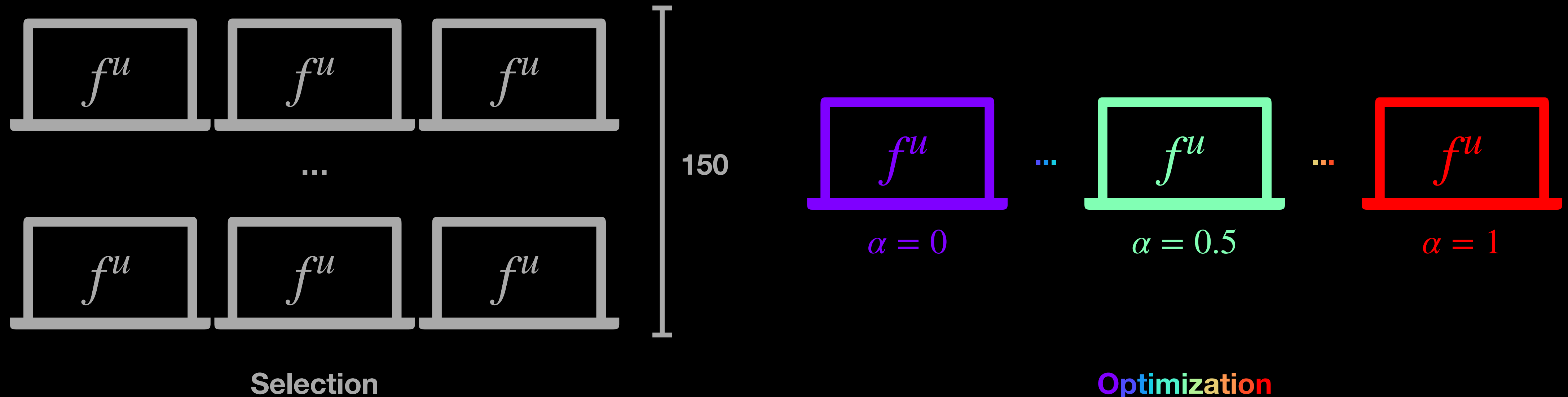
Updated Model Validation  
 $n = 2,500$

Evaluation  
Dataset  
 $n = 2,577$



$AUROC(f^o)$   
 $AUROC(f^u)$   
 $C^R(f^o, f^u)$

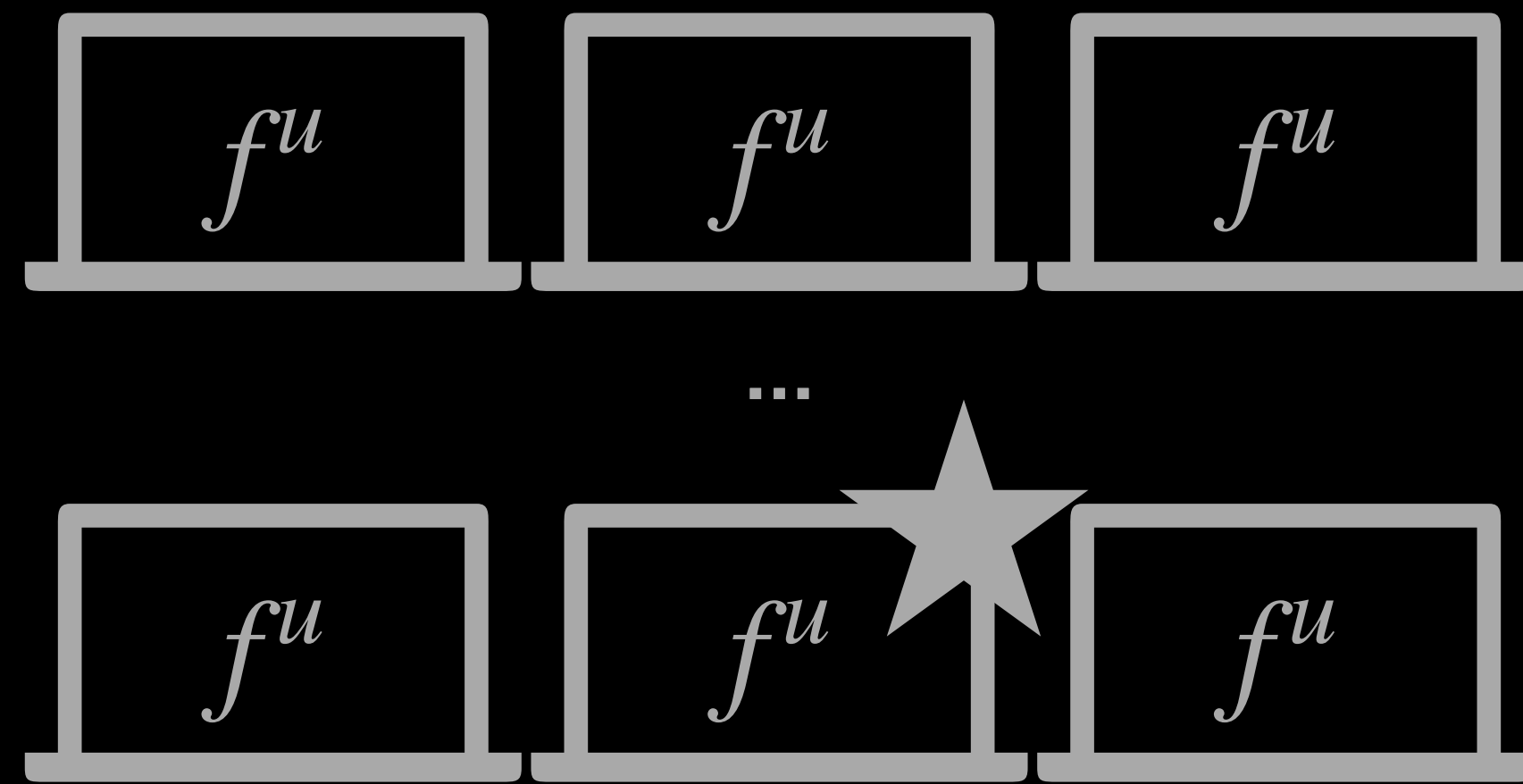
# Q2: Updated Models Selection vs. Optimization



150 “selection” models  
created through training with  
 $\mathcal{L}^{BCE}$  and randomly resampling  
the development dataset

11 “optimization” models  
created through training with  
weighted loss with  
 $\alpha = \{0, 0.1, \dots, 0.9, 1\}$

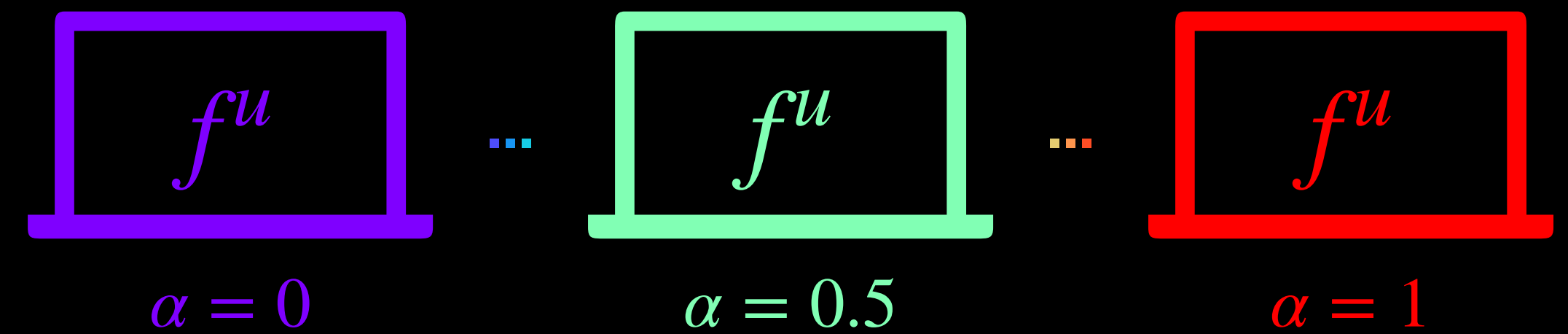
# Q2: Updated Models Selection vs. Optimization



Selection

Use a selection procedure to pick an updated model to use as a baseline

For example model with best validation AUROC

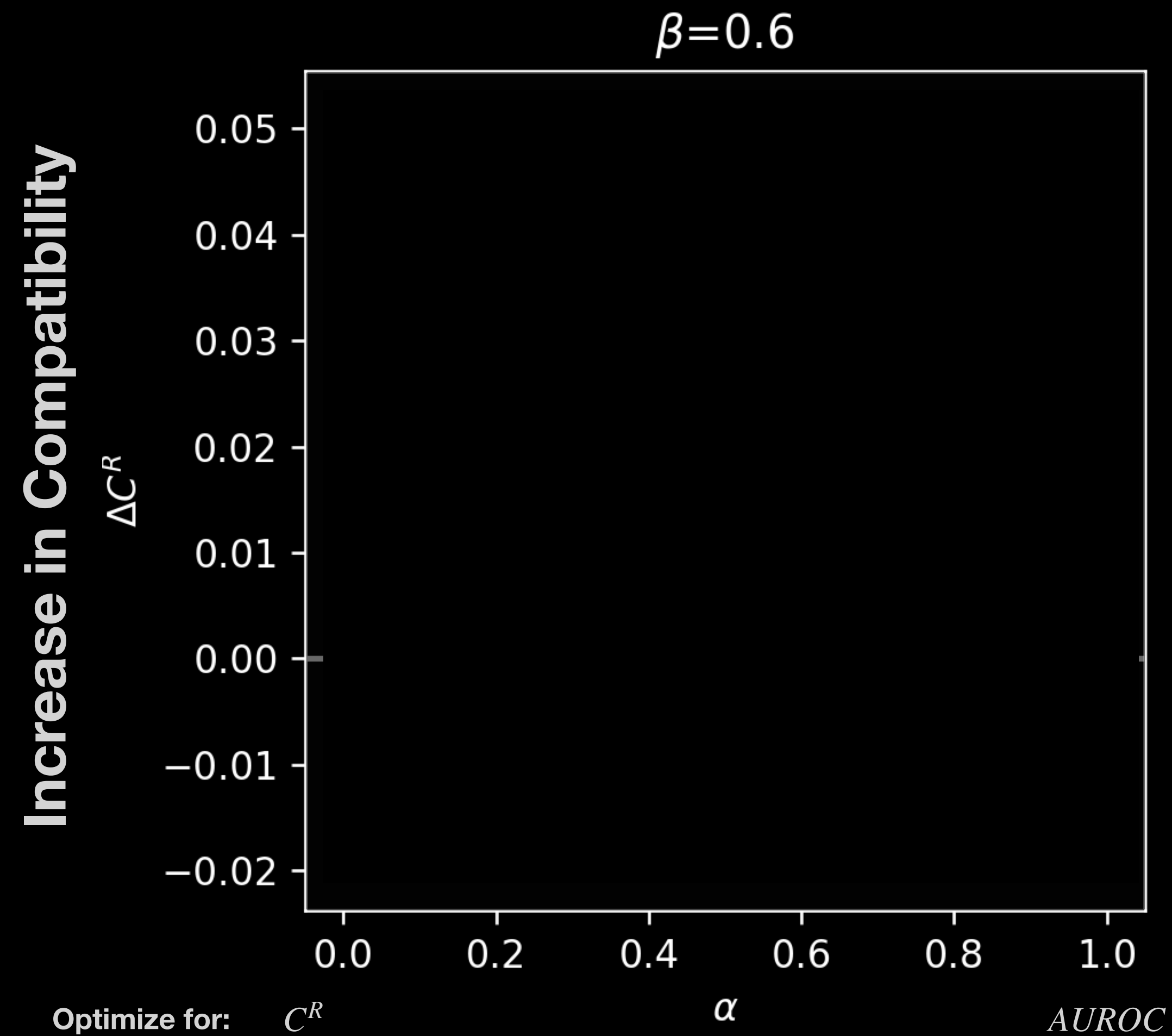


Optimization

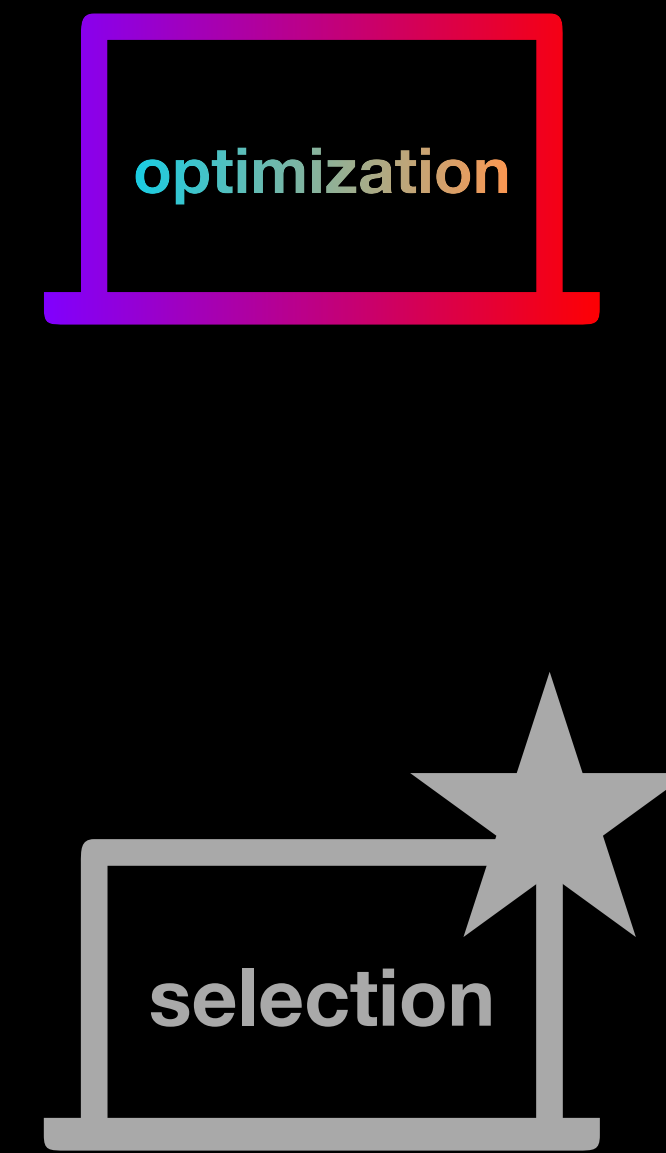
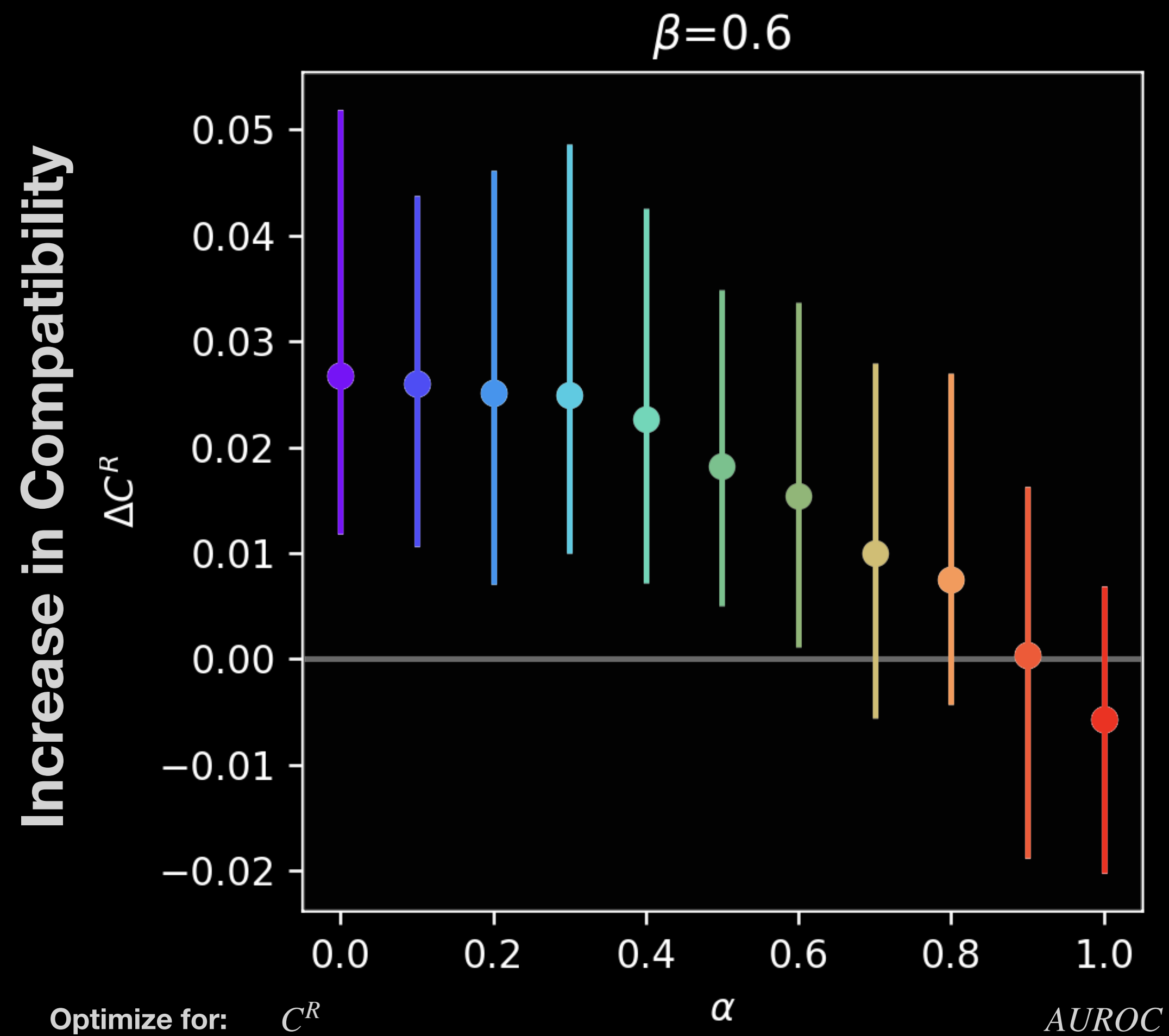
Examine difference in held out evaluation

$C^R$  and AUROC

# Q2: $C^R$ performance results

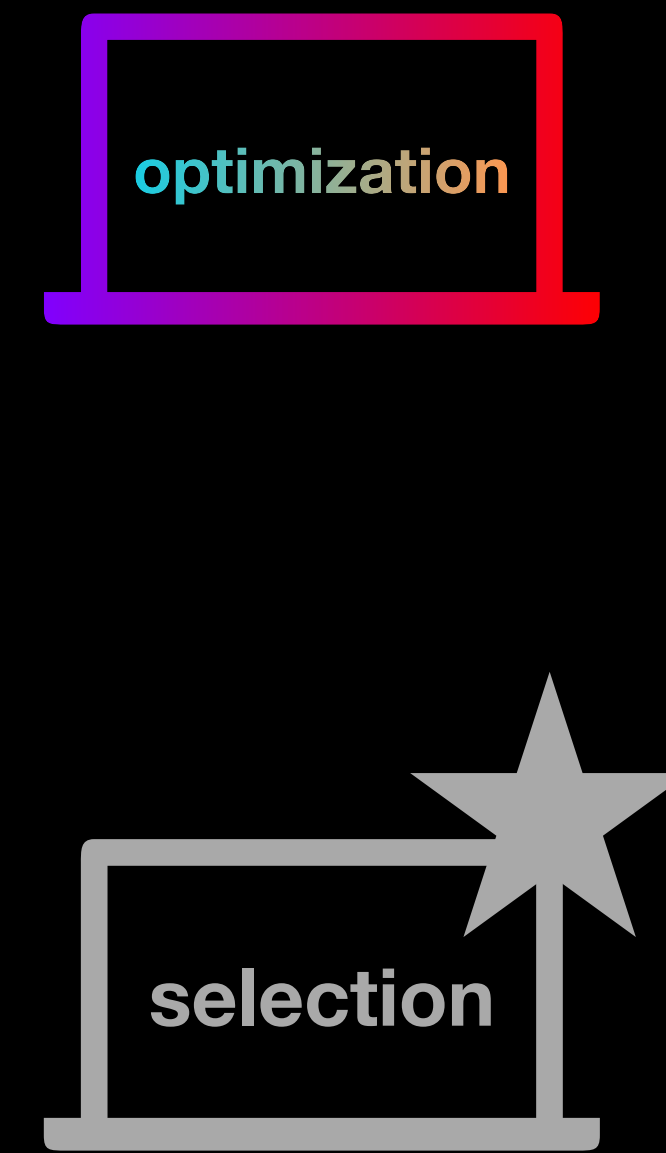
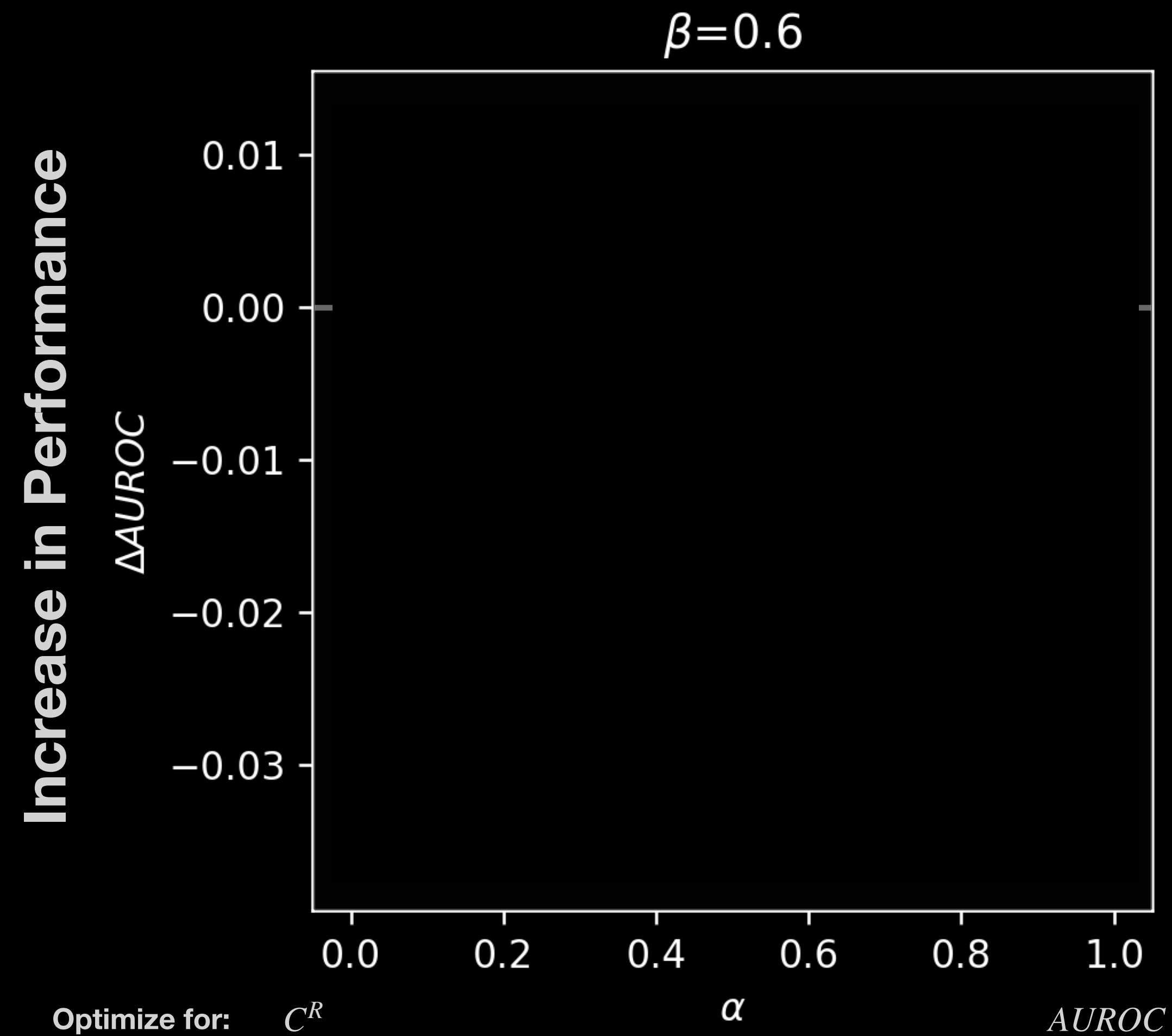


# Q2: $C^R$ performance results

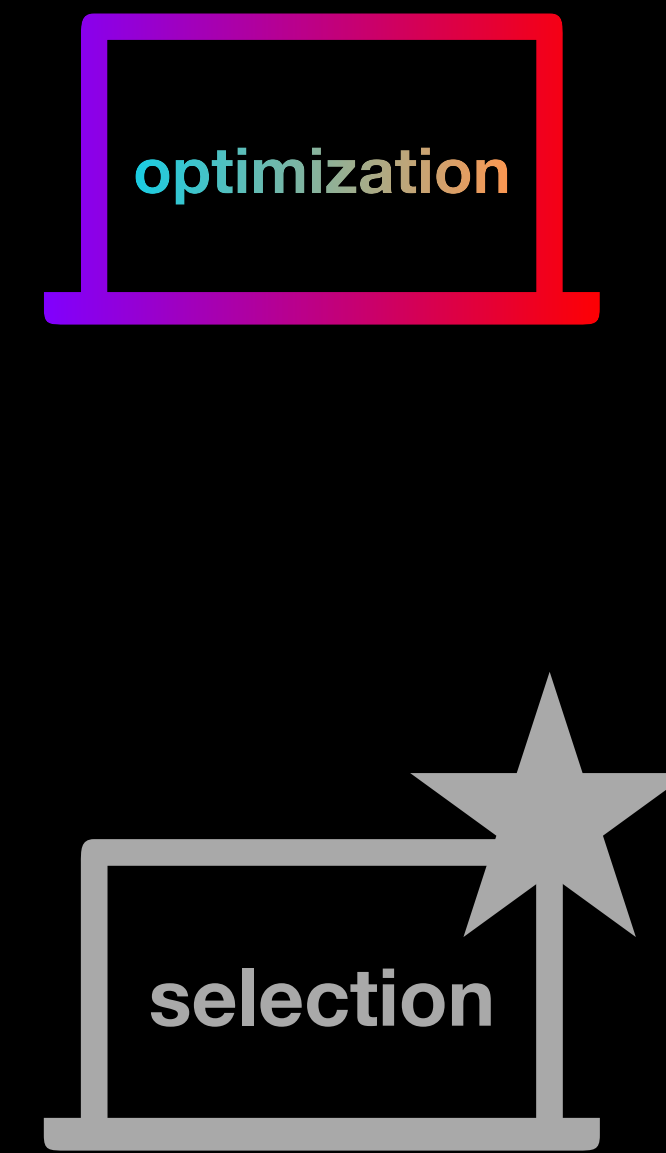
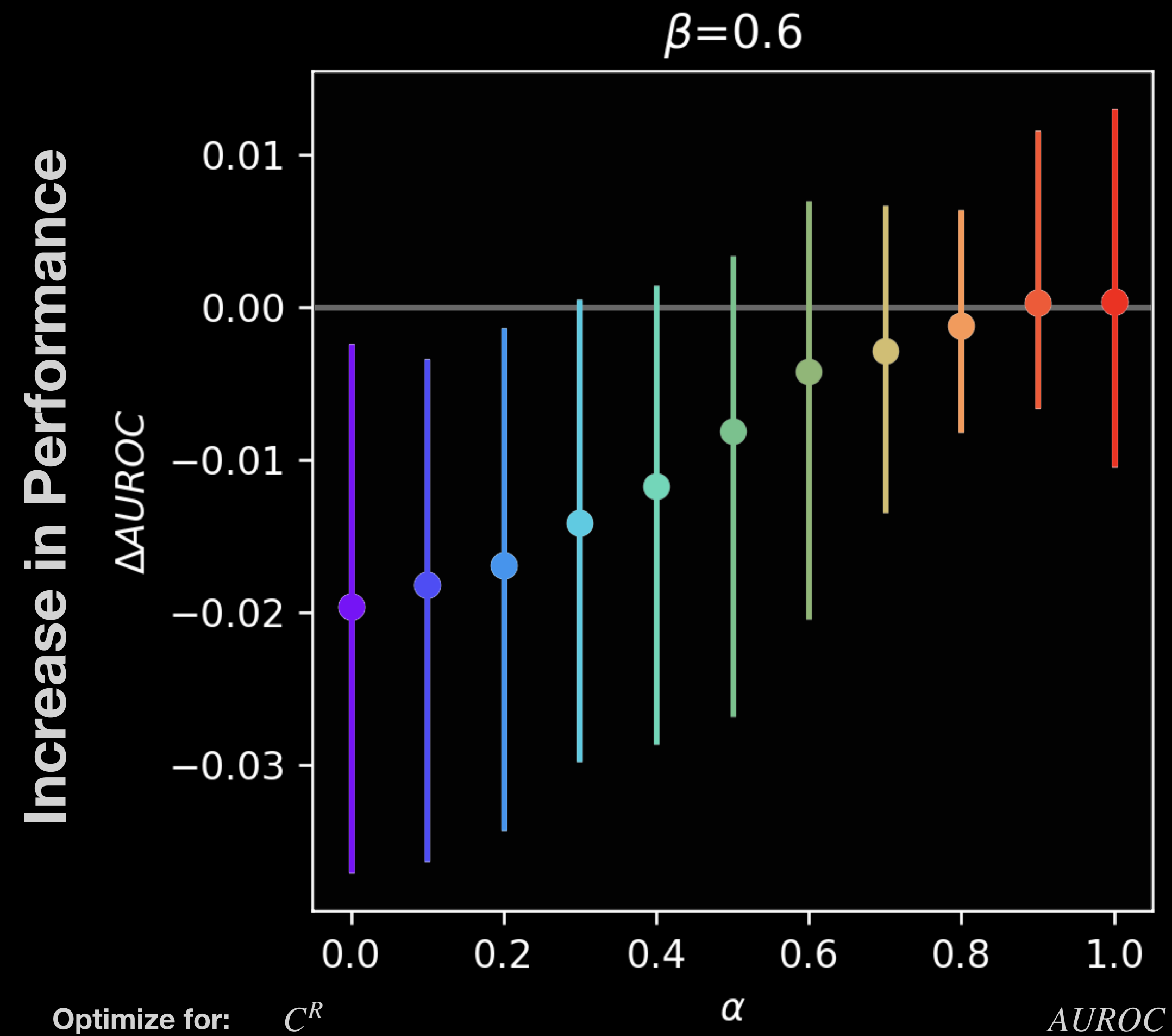




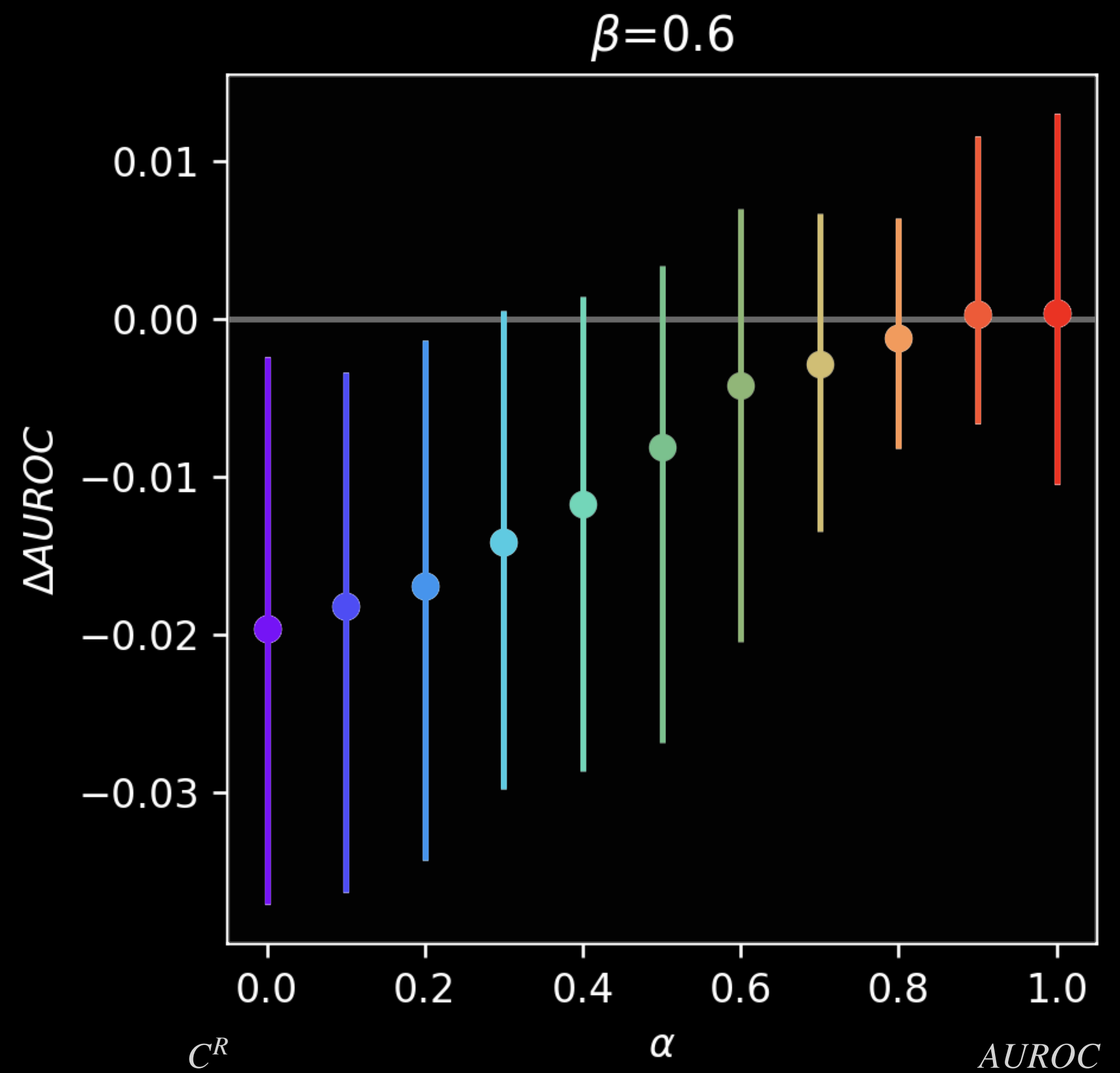
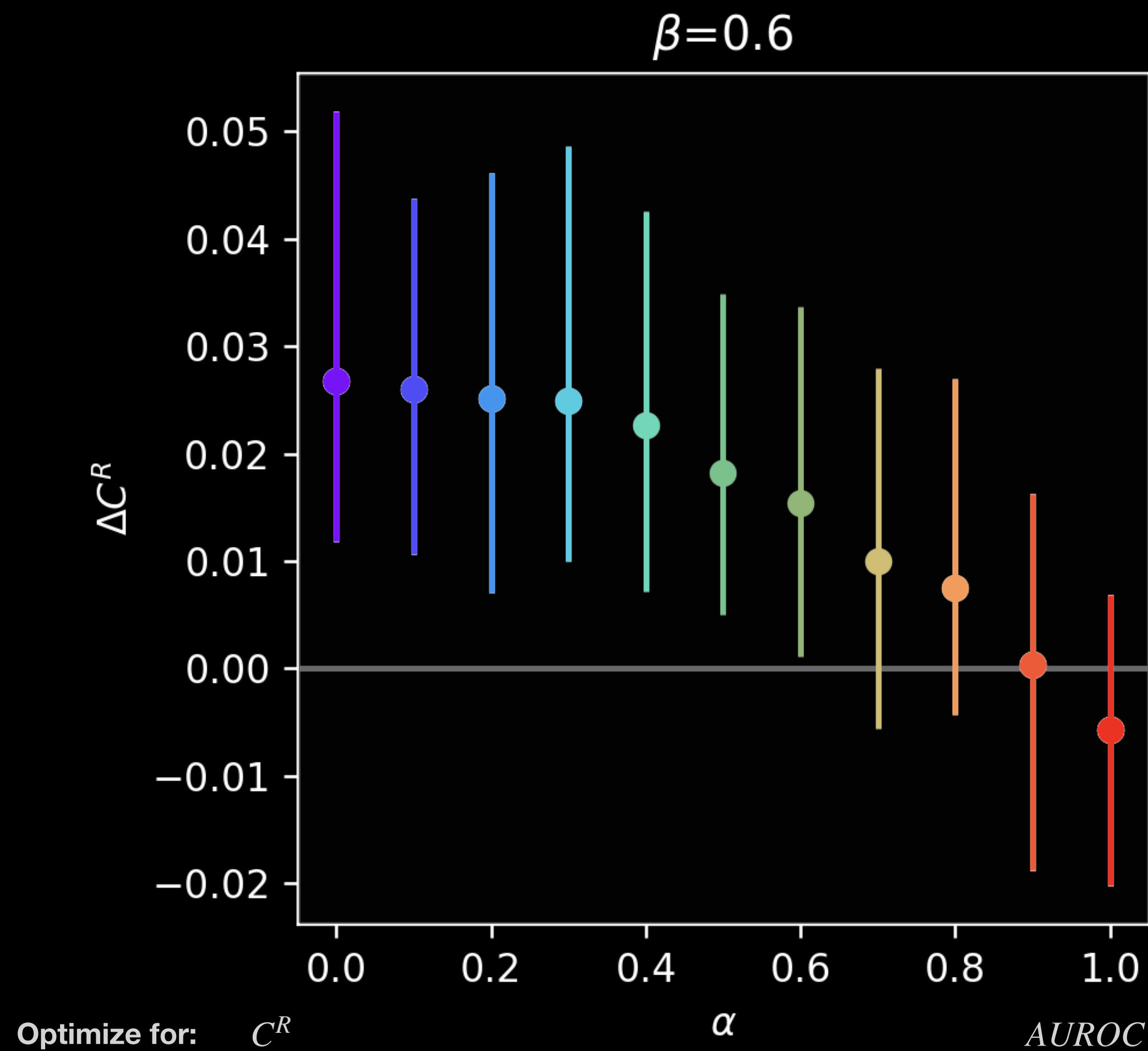
# Q2: *AUROC* performance results



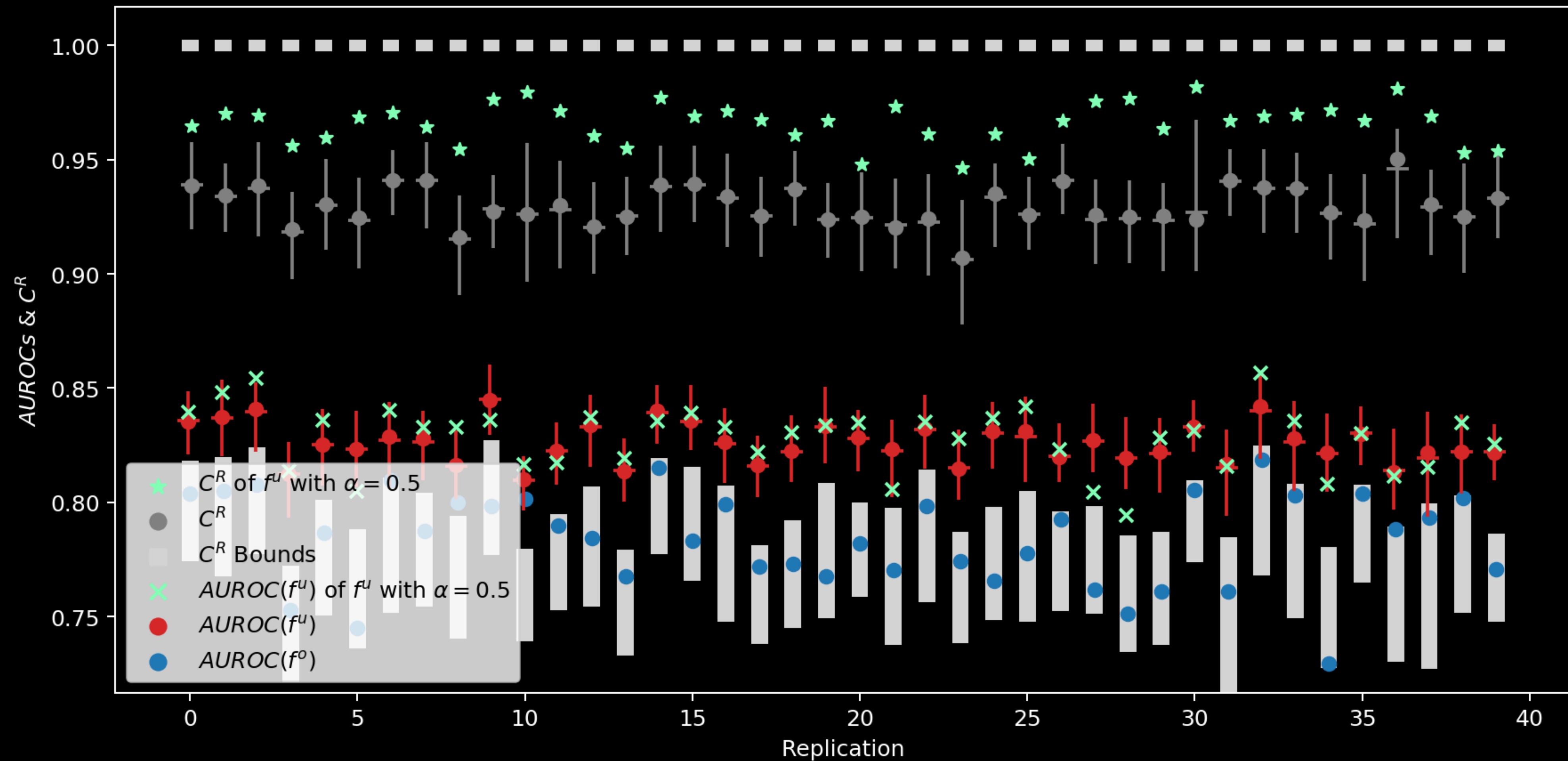
# Q2: *AUROC* performance results



# Q2: Performance results



# Q2: $\alpha=0.6$ yields promising updated models



# Summary of experiments

*Do we get  $C^R$  for free when we make updated models targeting AUROC?*

No.

*Can we make updated models with higher levels of  $C^R$ ?*

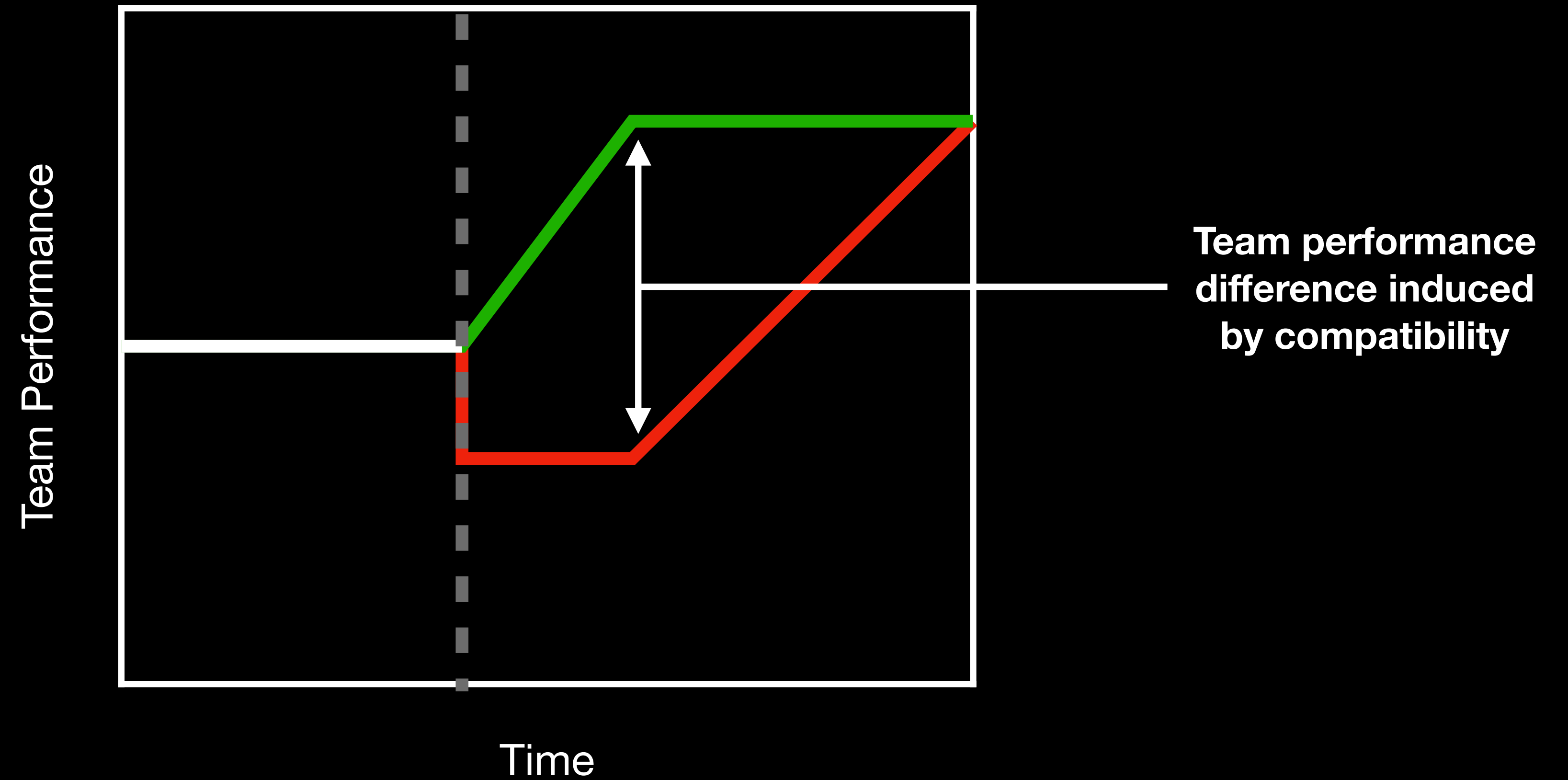
Yes, using our weighted loss function.

*Does that come at a cost in terms of AUROC?*

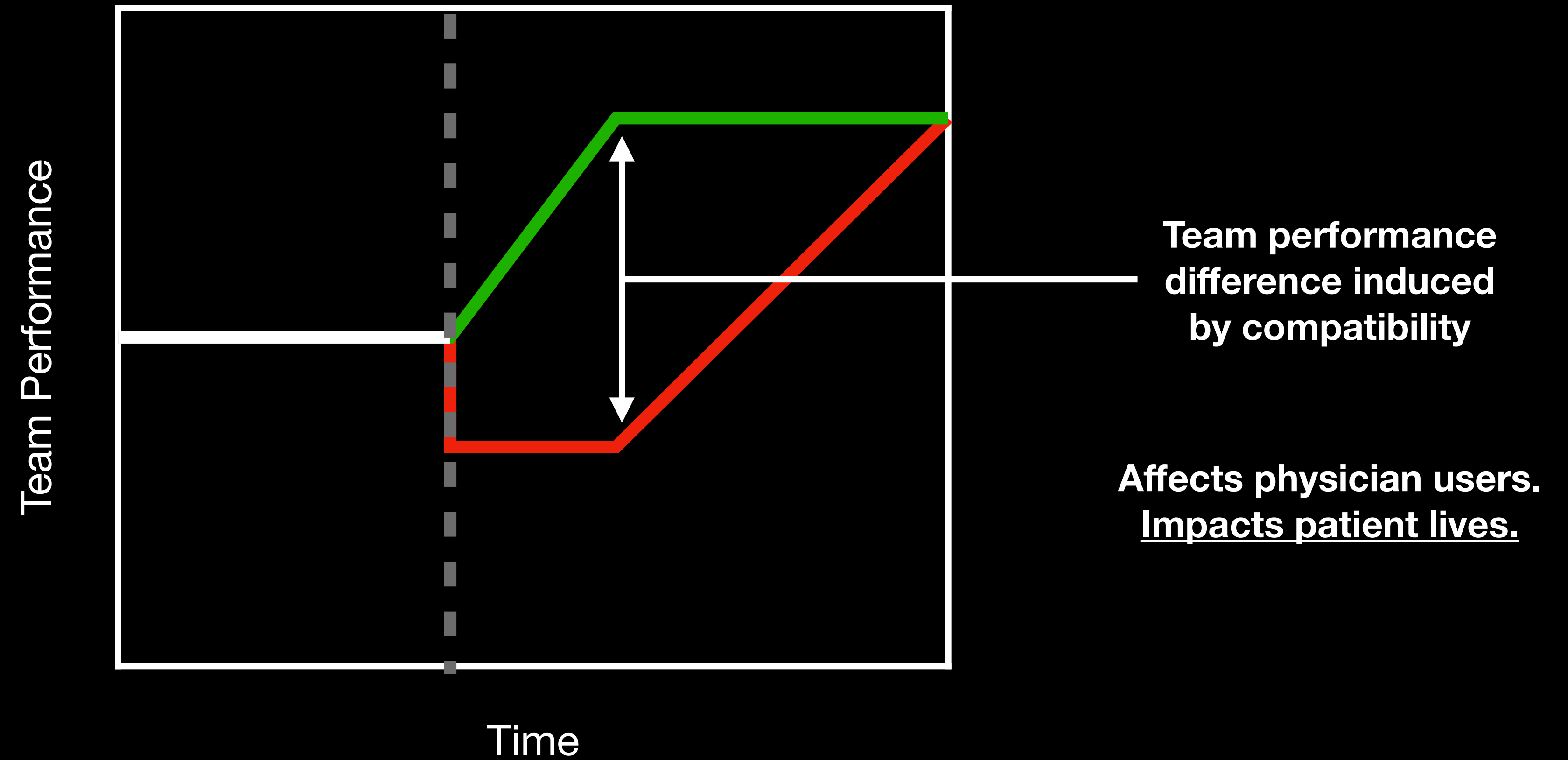
Sort of...

Why trade-off AUROC for  $C^R$ ?

# Why trade-off AUROC for $C^R$ ?

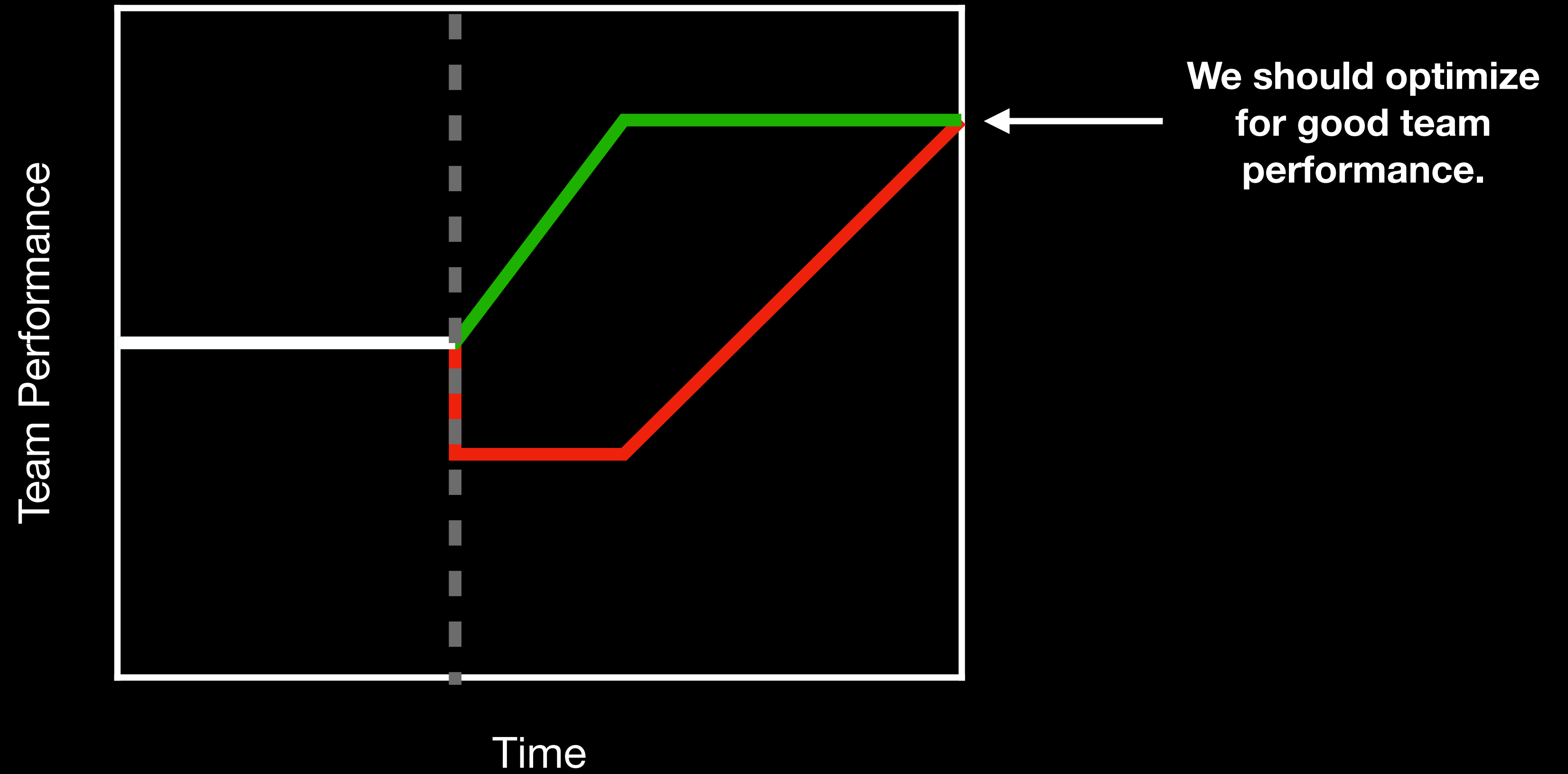


# Why trade-off AUROC for $C^R$ ?





# Why trade-off AUROC for $C^R$ ?



# $C^R$ is a new compatibility measure inspired by AUROC

Not threshold dependent:  $\uparrow$  clinical utility

Has direct relationship with AUROC

Can balance AUROC and  $C^R$

Using  $\widetilde{\mathcal{L}}^R \rightarrow \uparrow C^R$  &  $\uparrow$  AUROC

Real-world model updating case-study

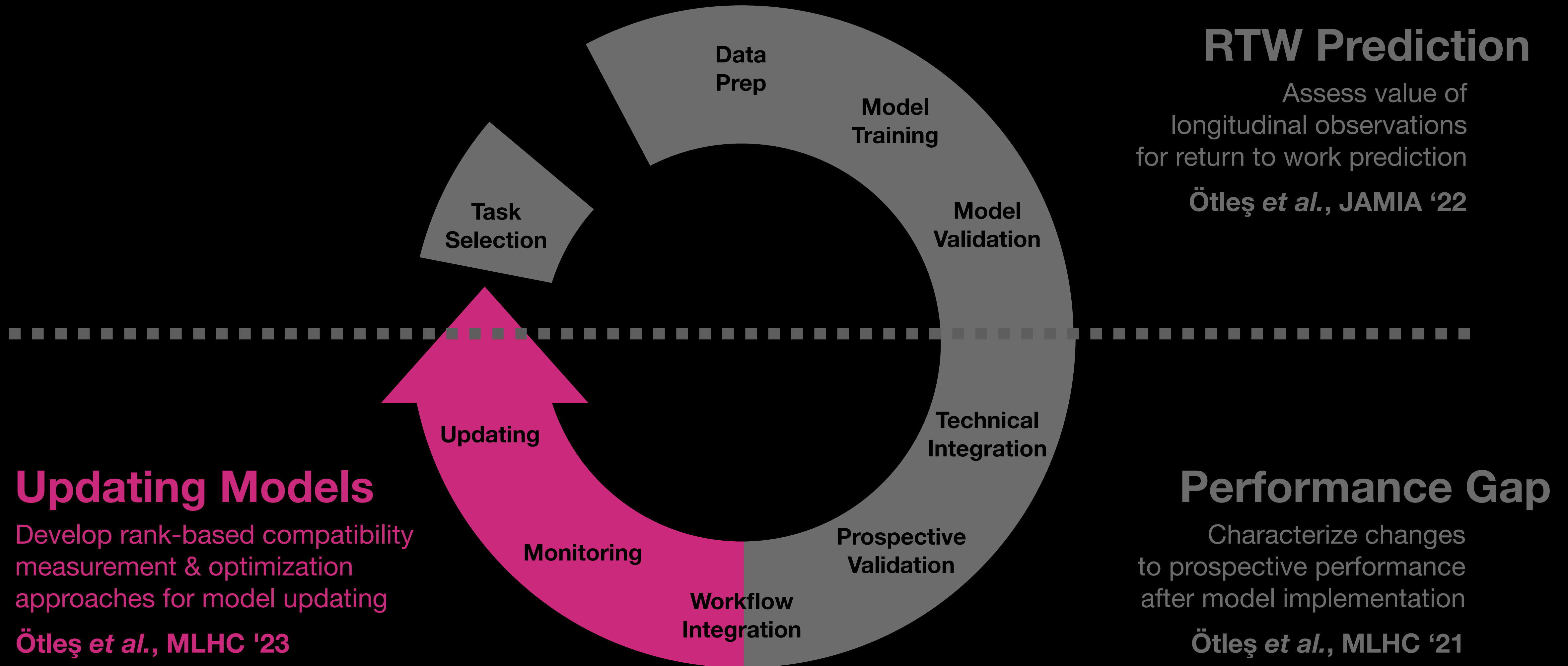
**My work enables model developers to find these desirable updated risk stratification models.**

Updated models may be assessed with our new rank-based compatibility measure.

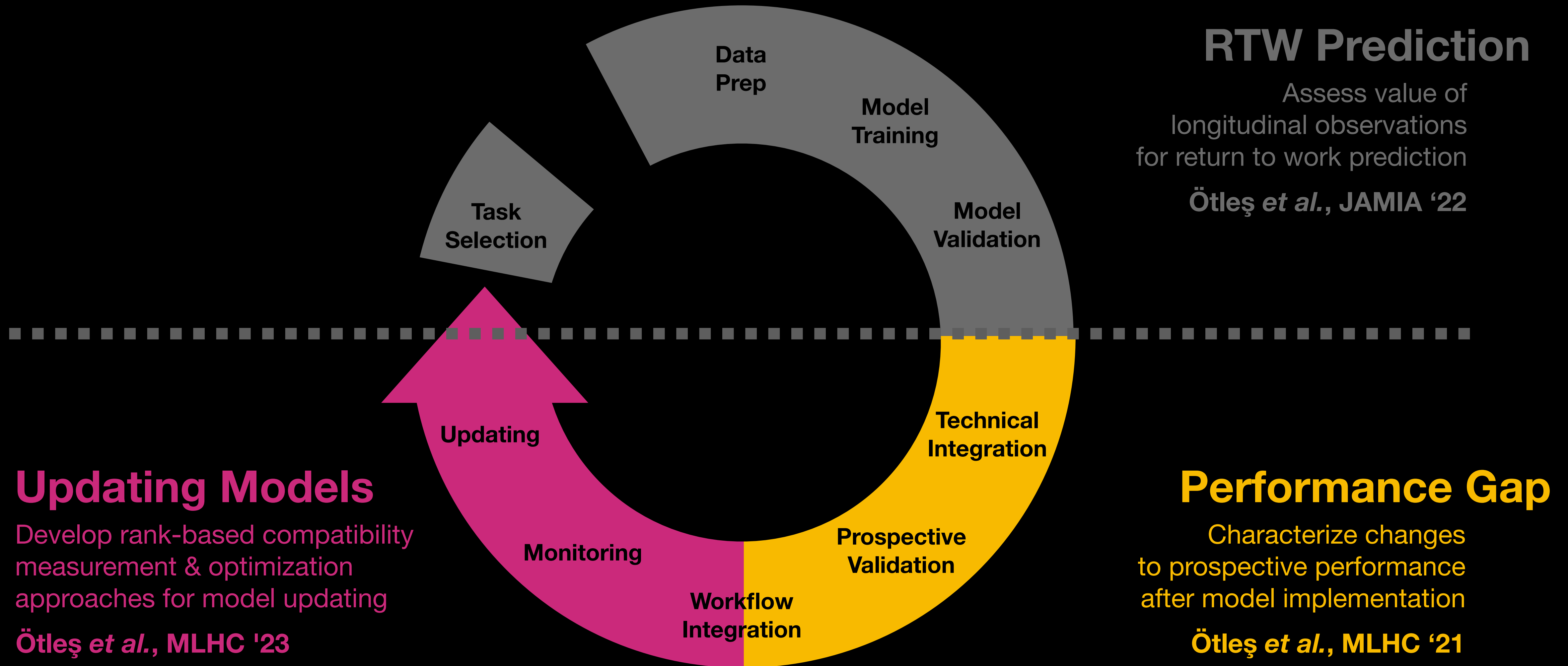
Model developers may use our weighted loss function to create model updates with improved  $C^R$ .

Additionally, we present a real-world model updating case-study.

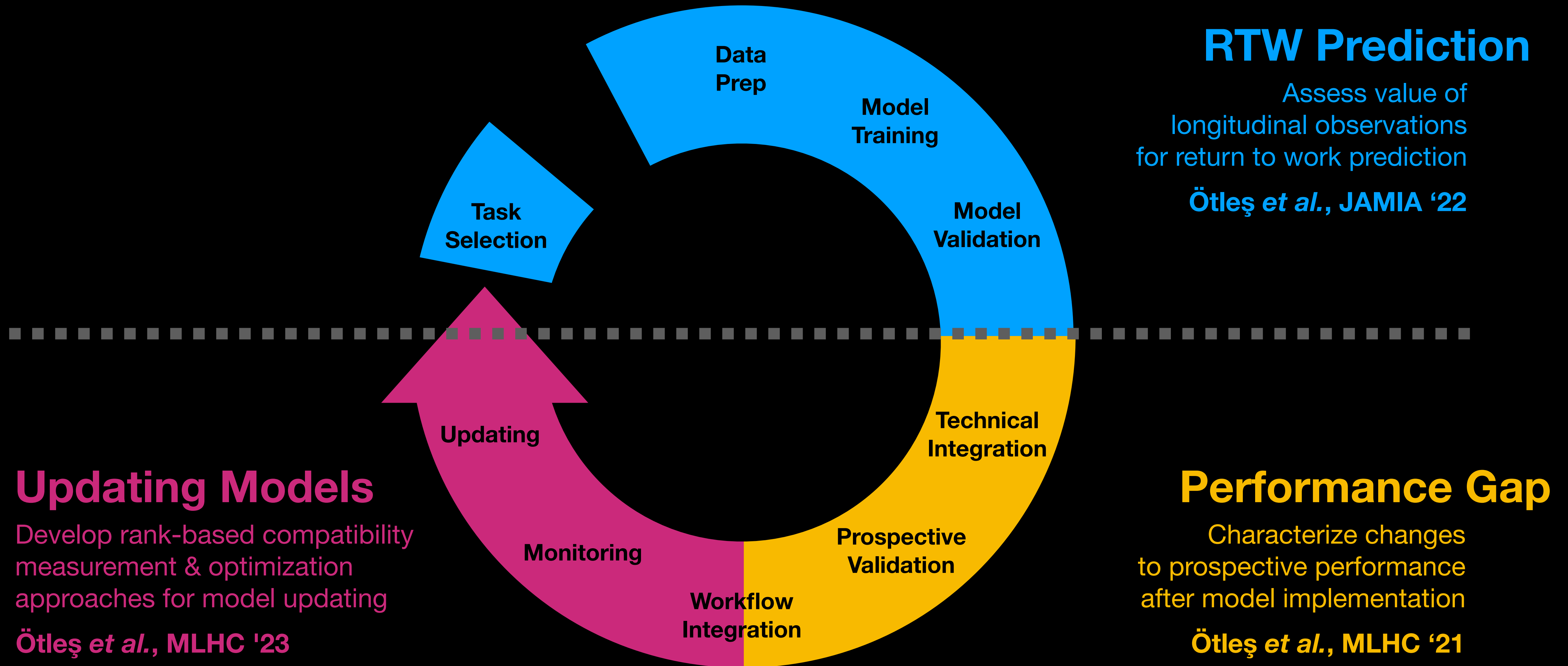
# The ML Lifecycle



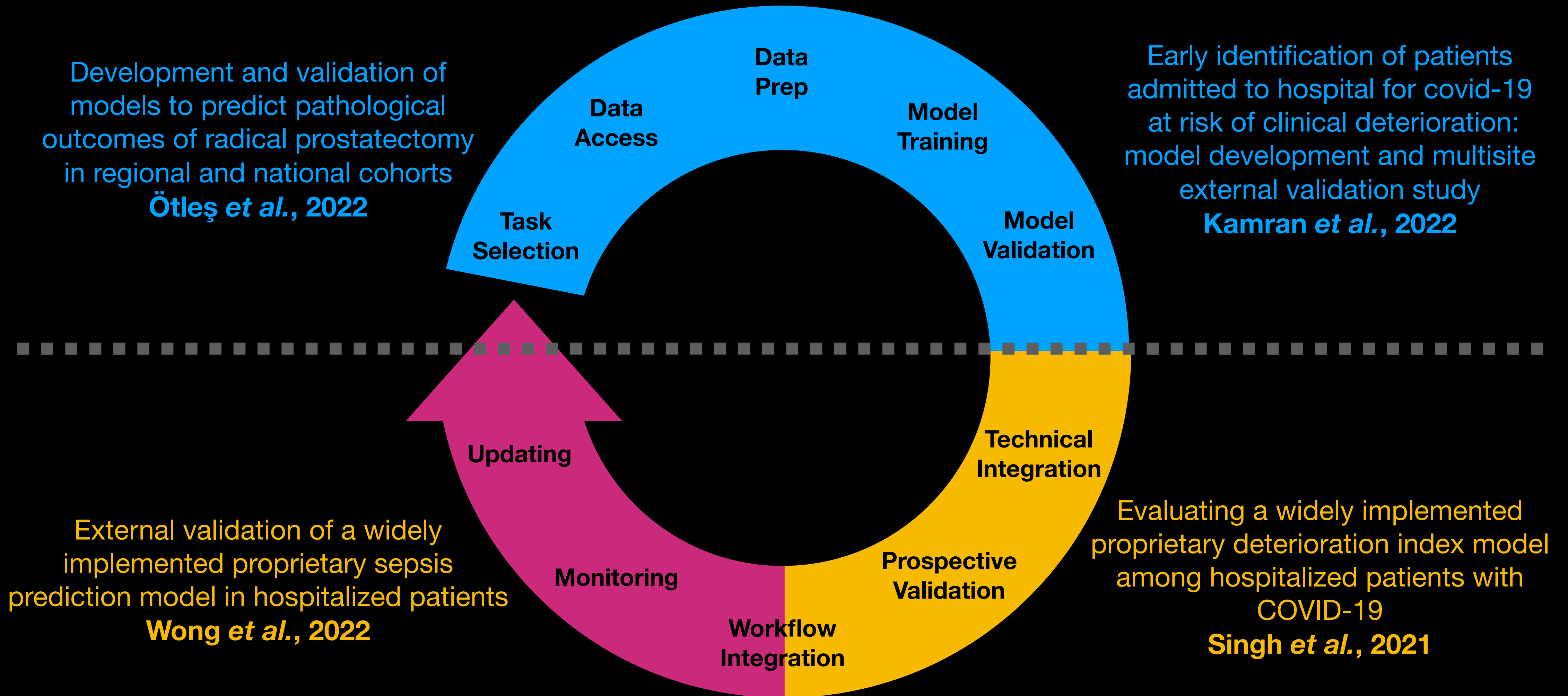
# The ML Lifecycle



# The ML Lifecycle

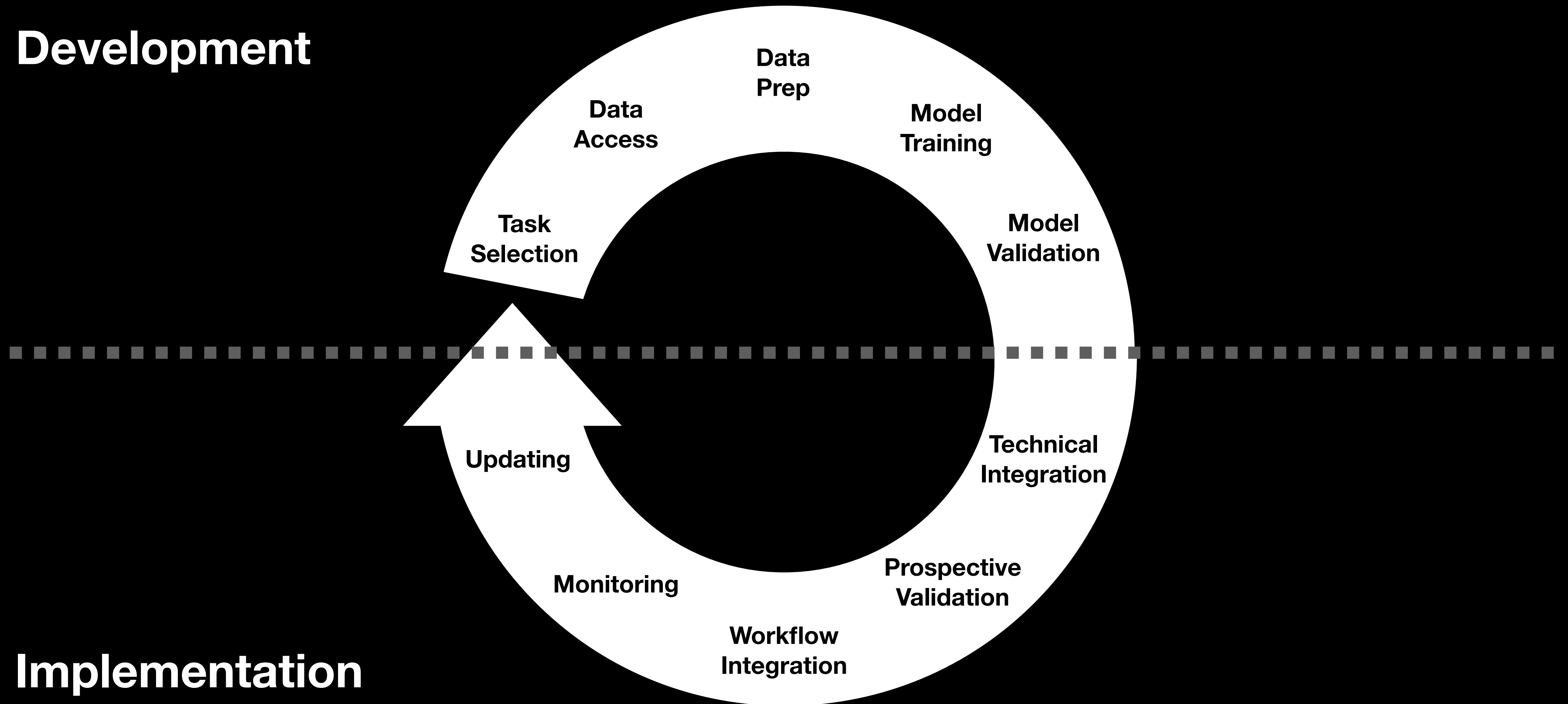


# Other projects



# Time plays a crucial role in the development and implementation of ML models used in healthcare.

**Development**



**Implementation**



# Future directions

Developing models less sensitive to infrastructure shifts

Updating

Personalized updating

Understanding user preferences for model updating

Updating models over time to ensure best team performance

Thesis: by focusing on the longitudinal nature of healthcare, we can improve the performance of predictive models during development & implementation.