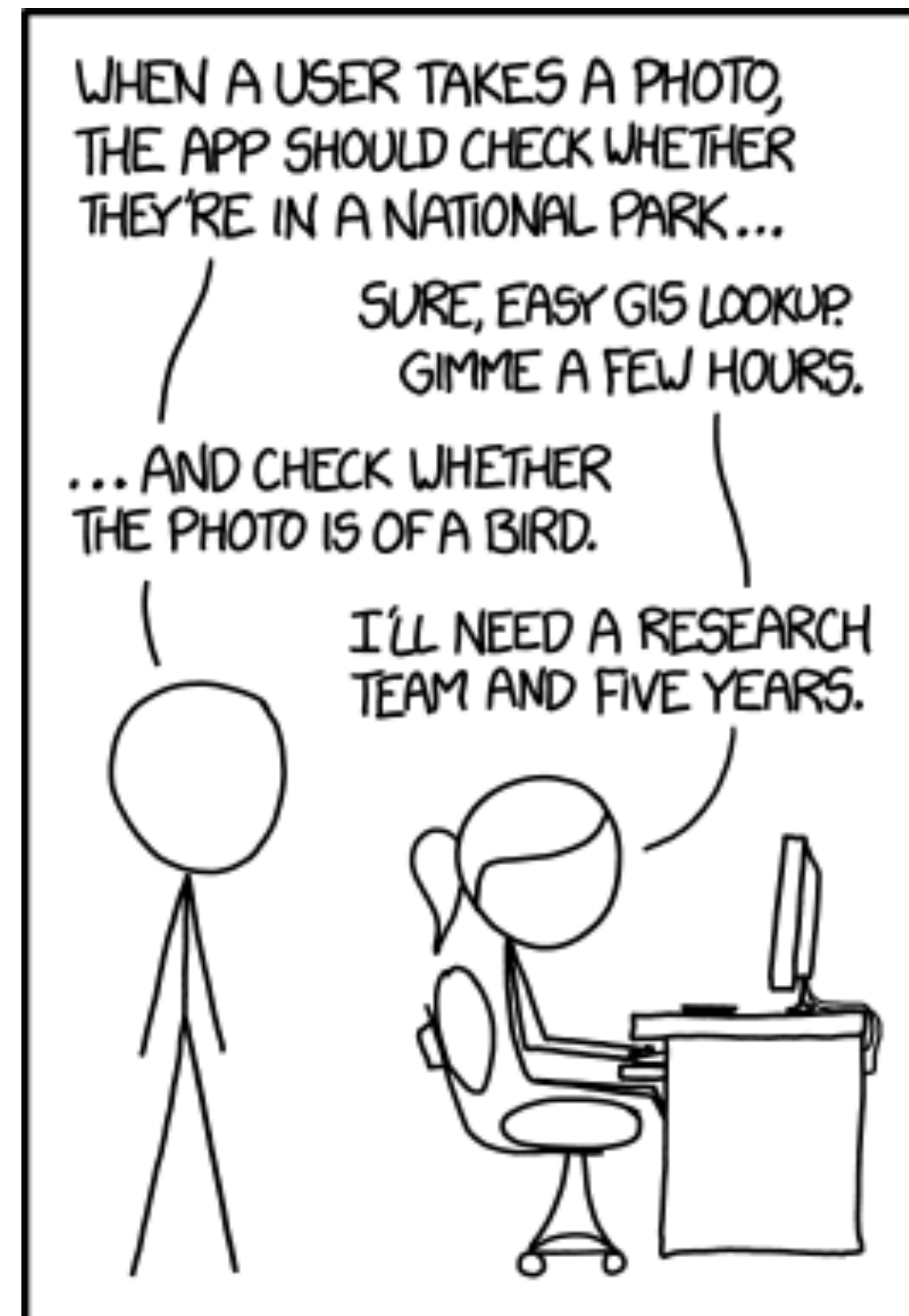


Setting up Machine Learning Projects

Josh Tobin, Sergey Karayev, Pieter Abbeel

Machine Learning Projects



IN CS, IT CAN BE HARD TO EXPLAIN
THE DIFFERENCE BETWEEN THE EASY
AND THE VIRTUALLY IMPOSSIBLE.

Machine Learning Projects

85% of AI projects fail¹

¹ Pactera Technologies

Why do so many projects fail?

- ML is still research - you shouldn't aim for 100% success rate
- But, many are doomed to fail:
- Technically infeasible or poorly scoped
- Never make the leap to production
- Unclear success criteria
- Poor team management

Module overview

-  **Lifecycle**
 - How to think about all of the activities in an ML project
-  **Prioritizing projects**
 - Assessing the feasibility and impact of your projects
-  **Archetypes**
 - The main categories of ML projects, and the implications for project management
-  **Metrics**
 - How to pick a single number to optimize
-  **Baselines**
 - How to know if your model is performing well

Running case study - pose estimation



$$(x, y, z)$$

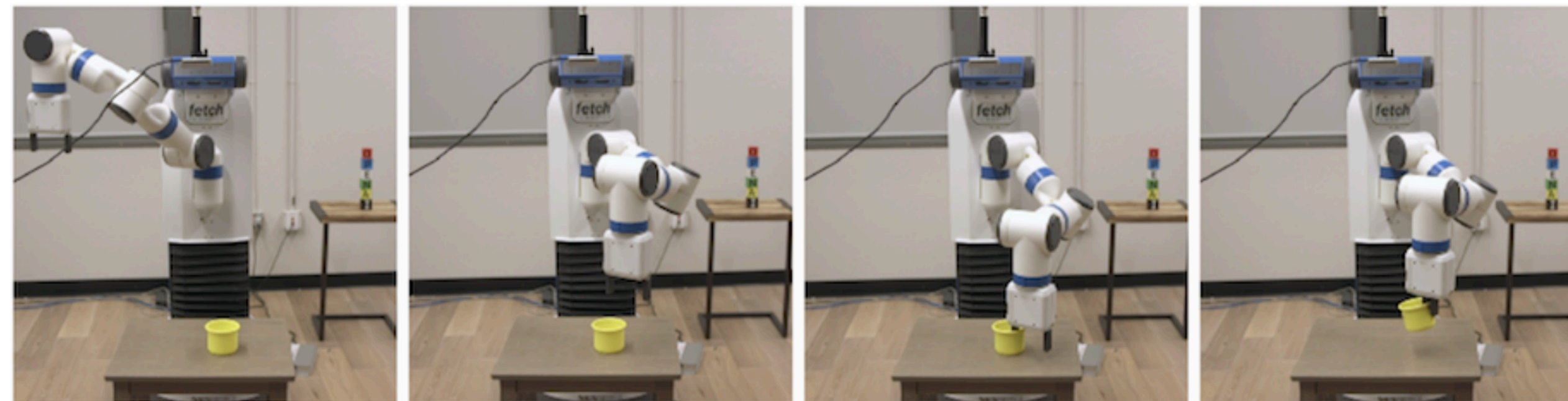
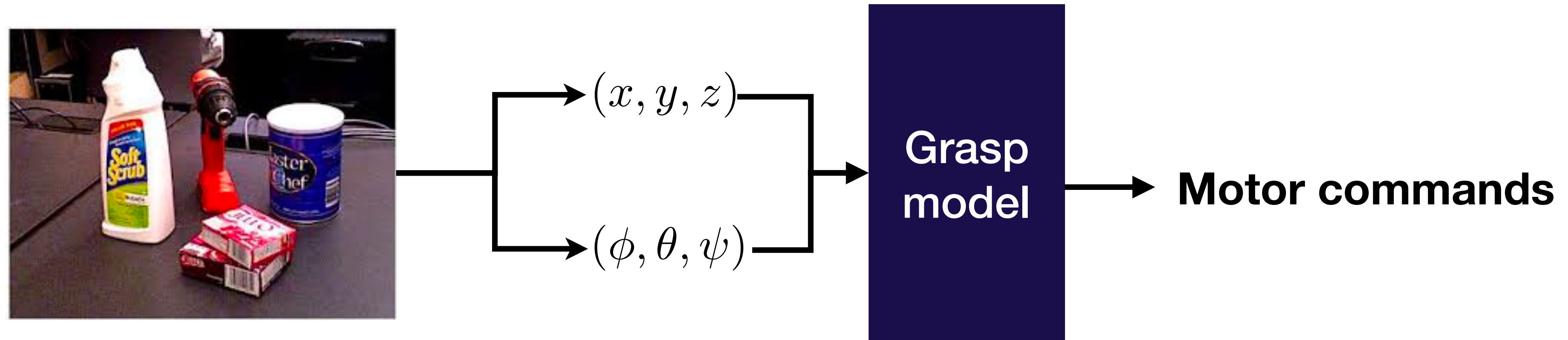
Position (L2 loss)

$$(\phi, \theta, \psi)$$

Orientation (L2 loss)

Xiang, Yu, et al. "PoseCNN: A Convolutional Neural Network for 6D Object Pose Estimation in Cluttered Scenes." arXiv preprint arXiv:1711.00199 (2017).

Full Stack Robotics works on grasping



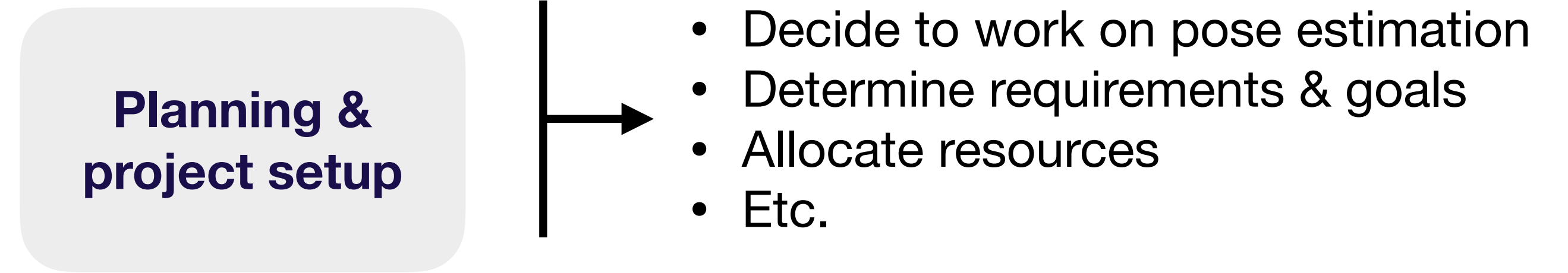
Module overview

-  **Lifecycle**
 - **How to think about all of the activities in an ML project**
-  **Prioritizing projects**
 - Assessing the feasibility and impact of your projects
-  **Archetypes**
 - The main categories of ML projects, and the implications for project management
-  **Metrics**
 - How to pick a single number to optimize
-  **Baselines**
 - How to know if your model is performing well

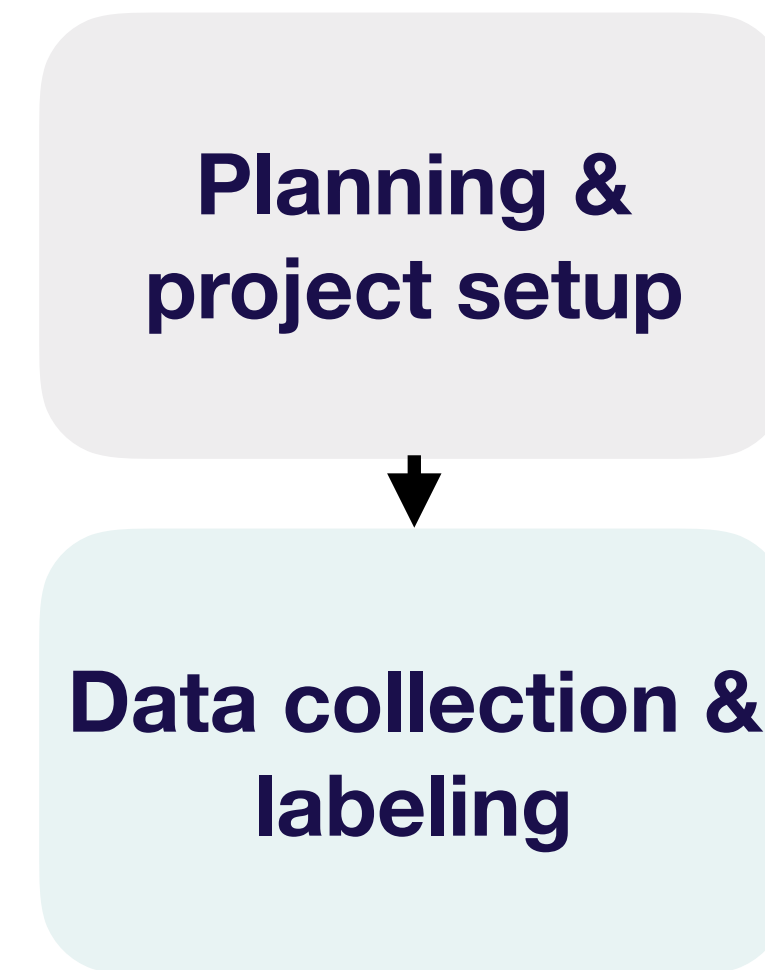
Lifecycle of a ML project

**Planning &
project setup**

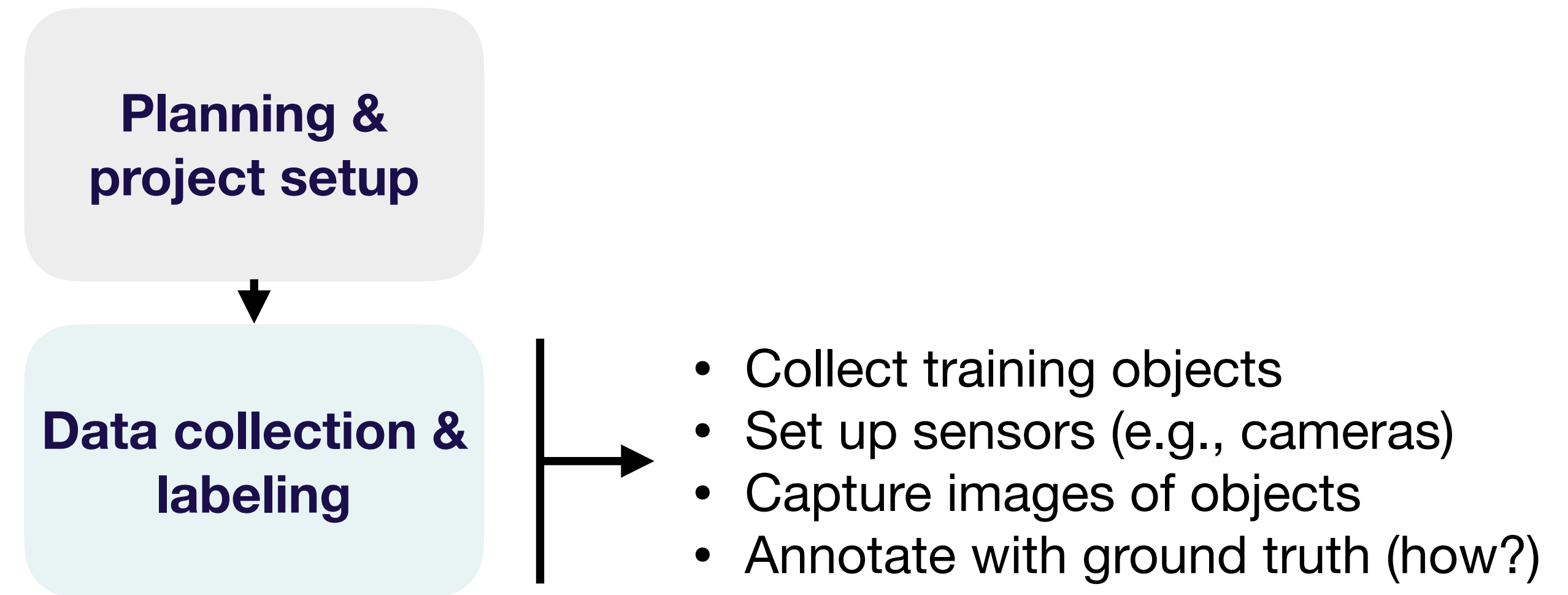
Lifecycle of a ML project



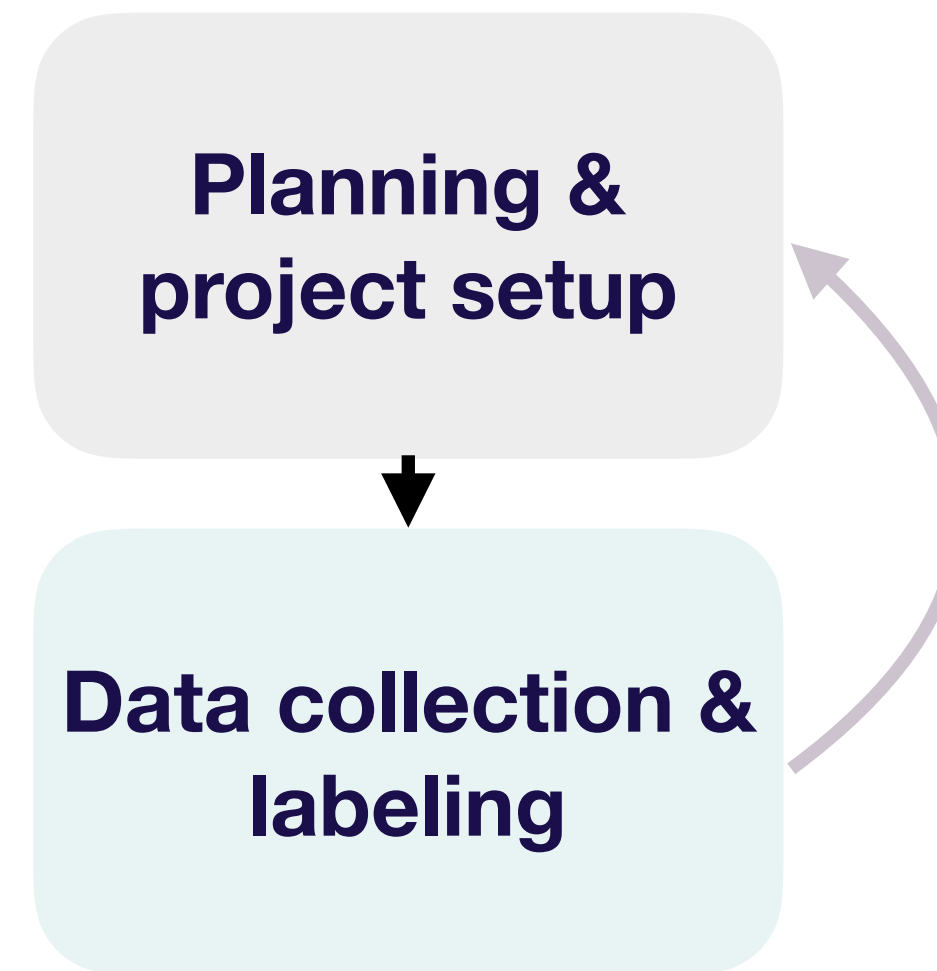
Lifecycle of a ML project



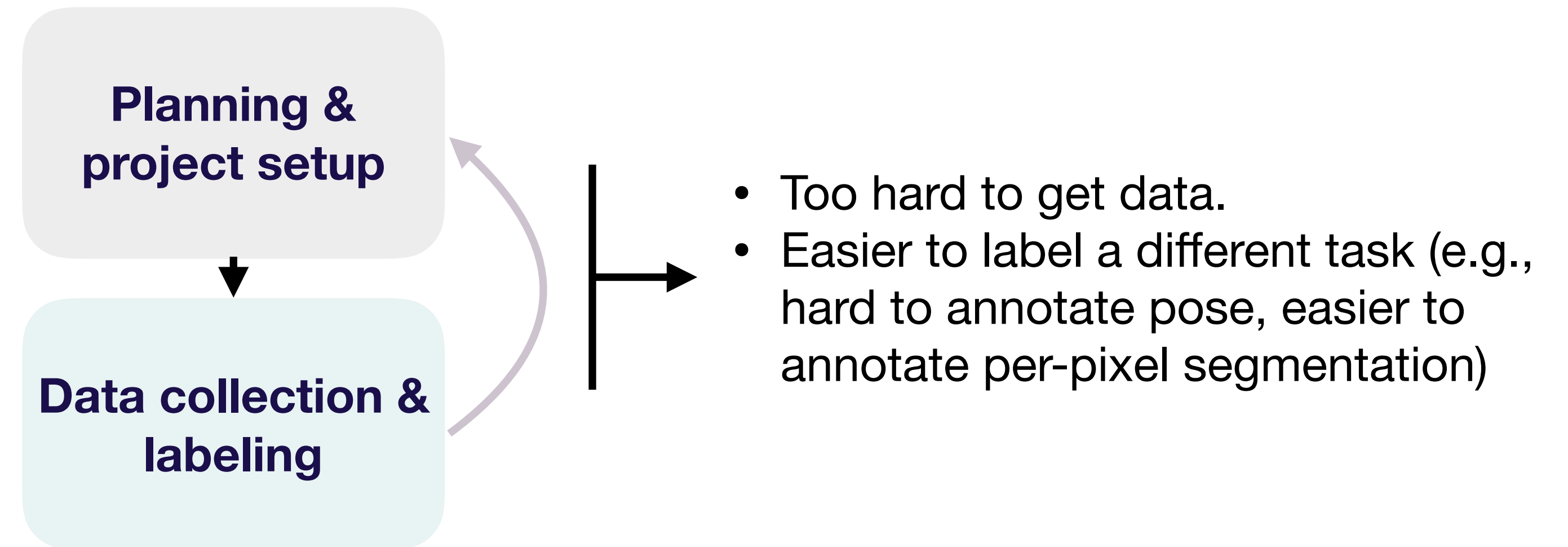
Lifecycle of a ML project



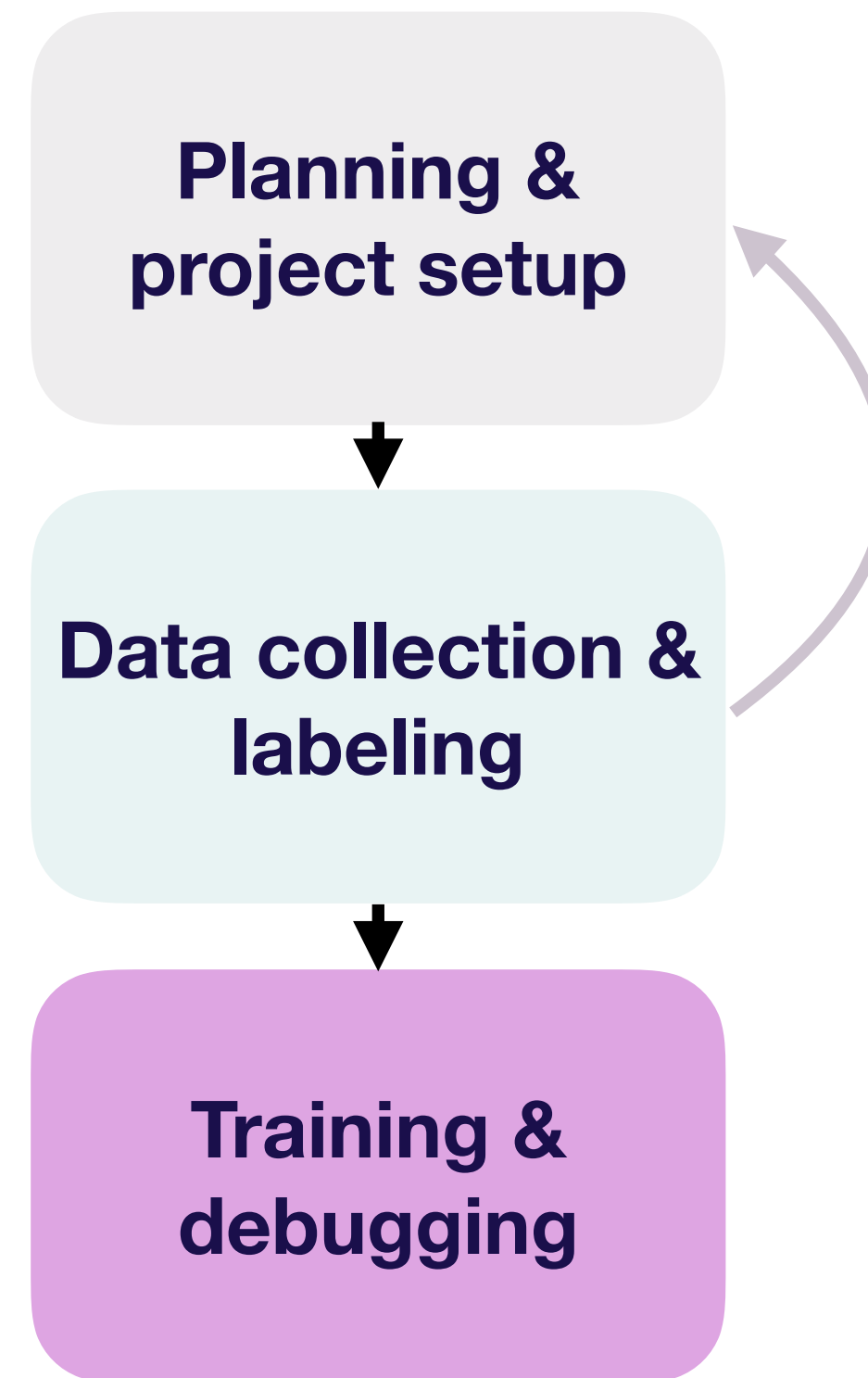
Lifecycle of a ML project



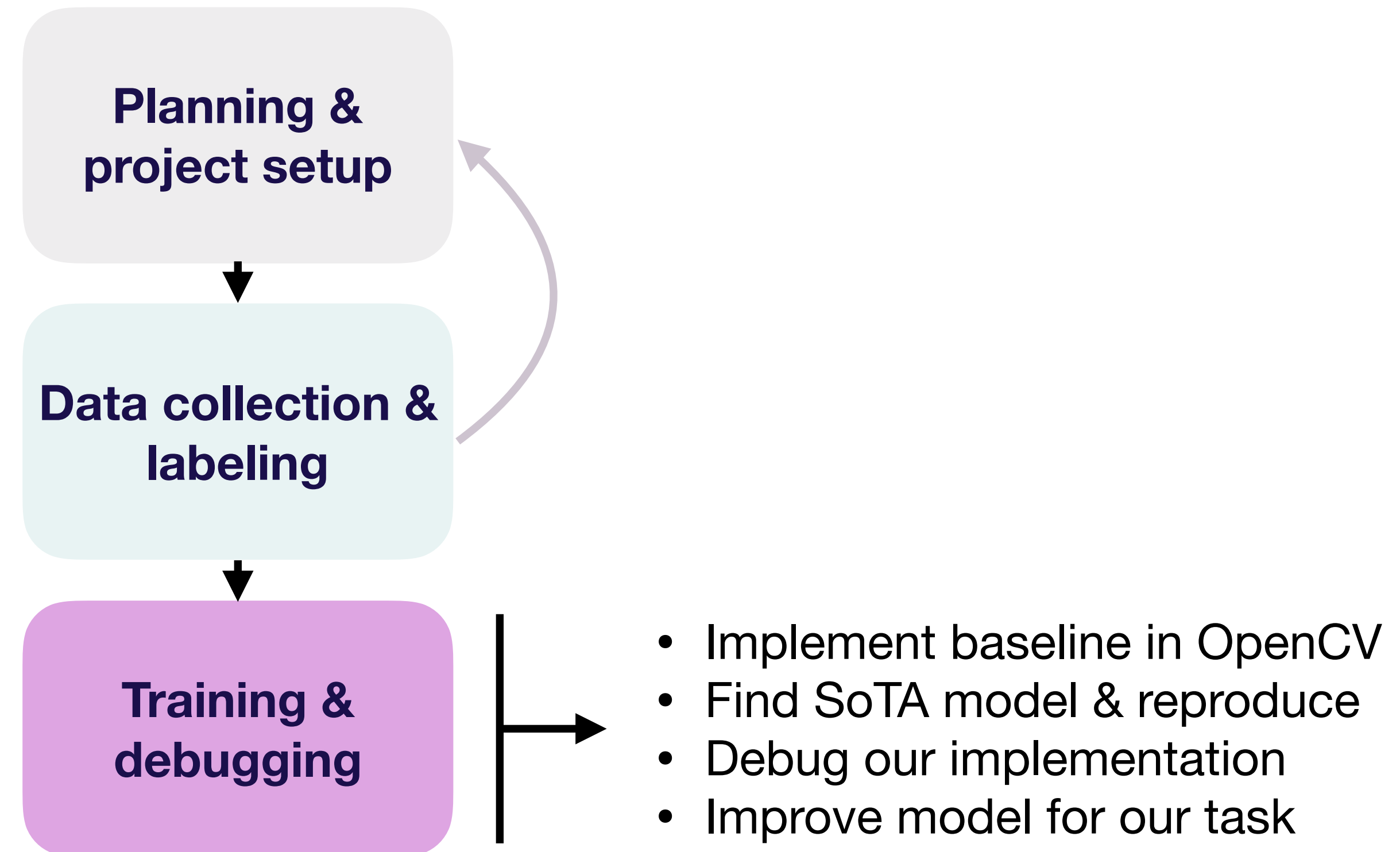
Lifecycle of a ML project



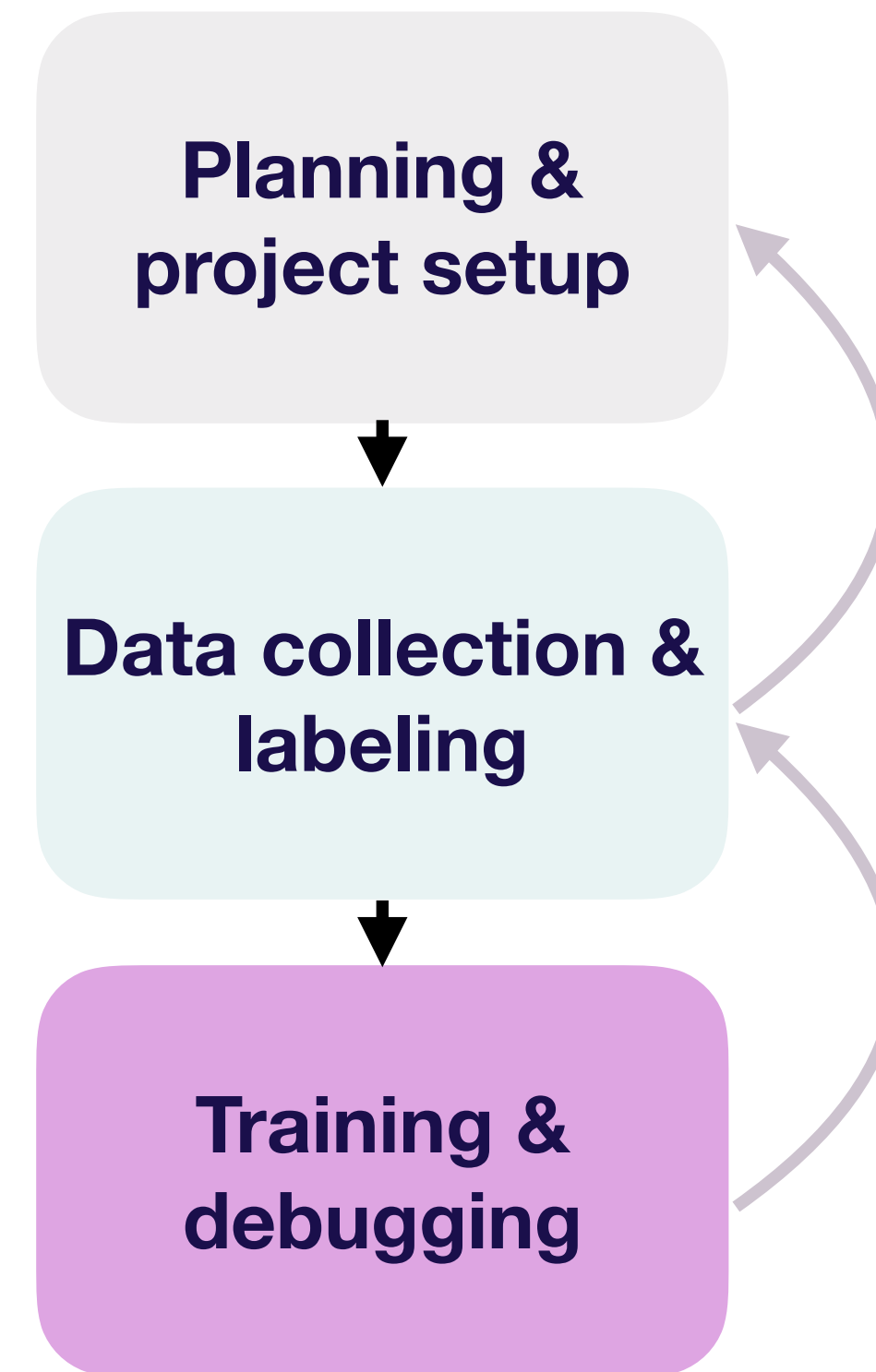
Lifecycle of a ML project



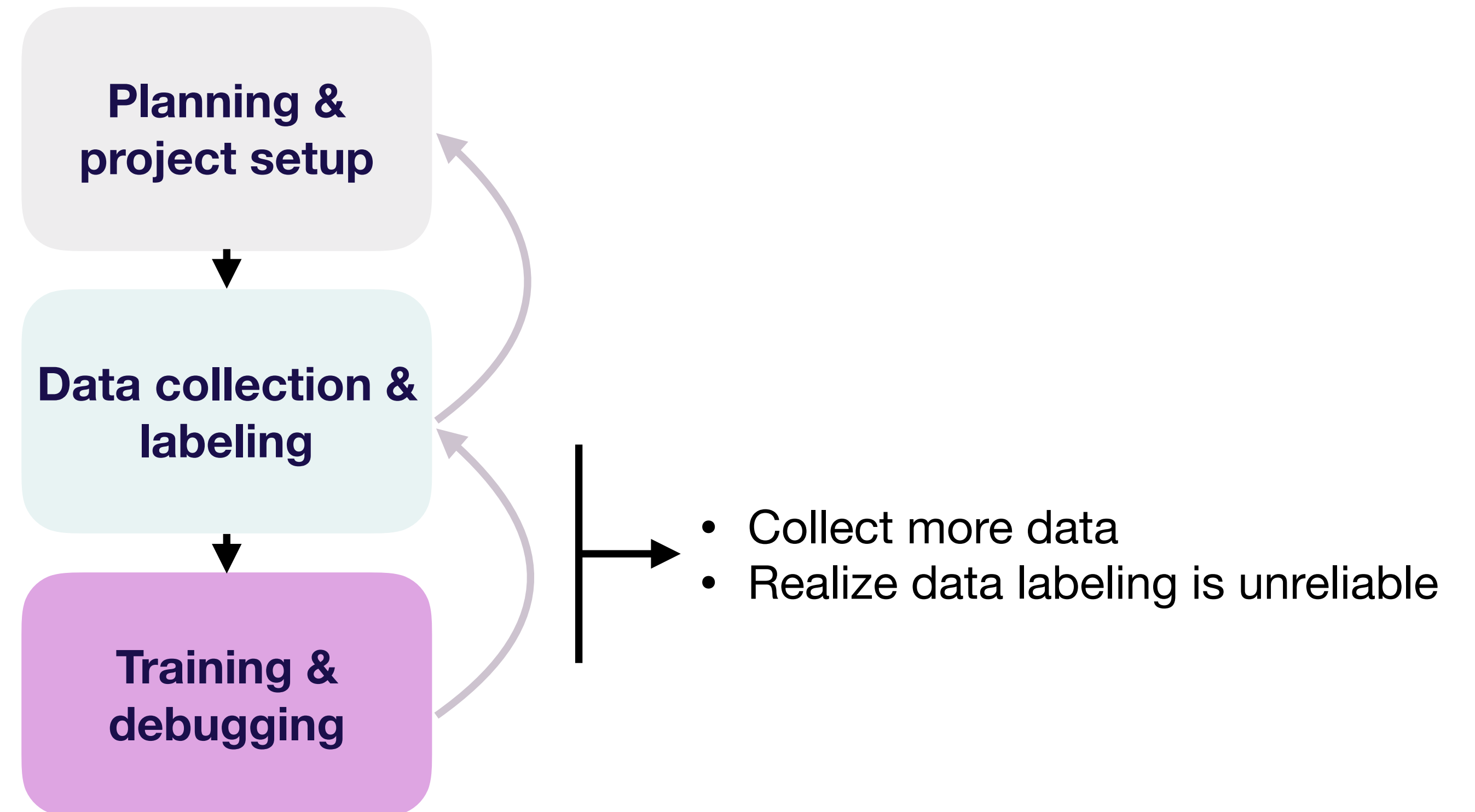
Lifecycle of a ML project



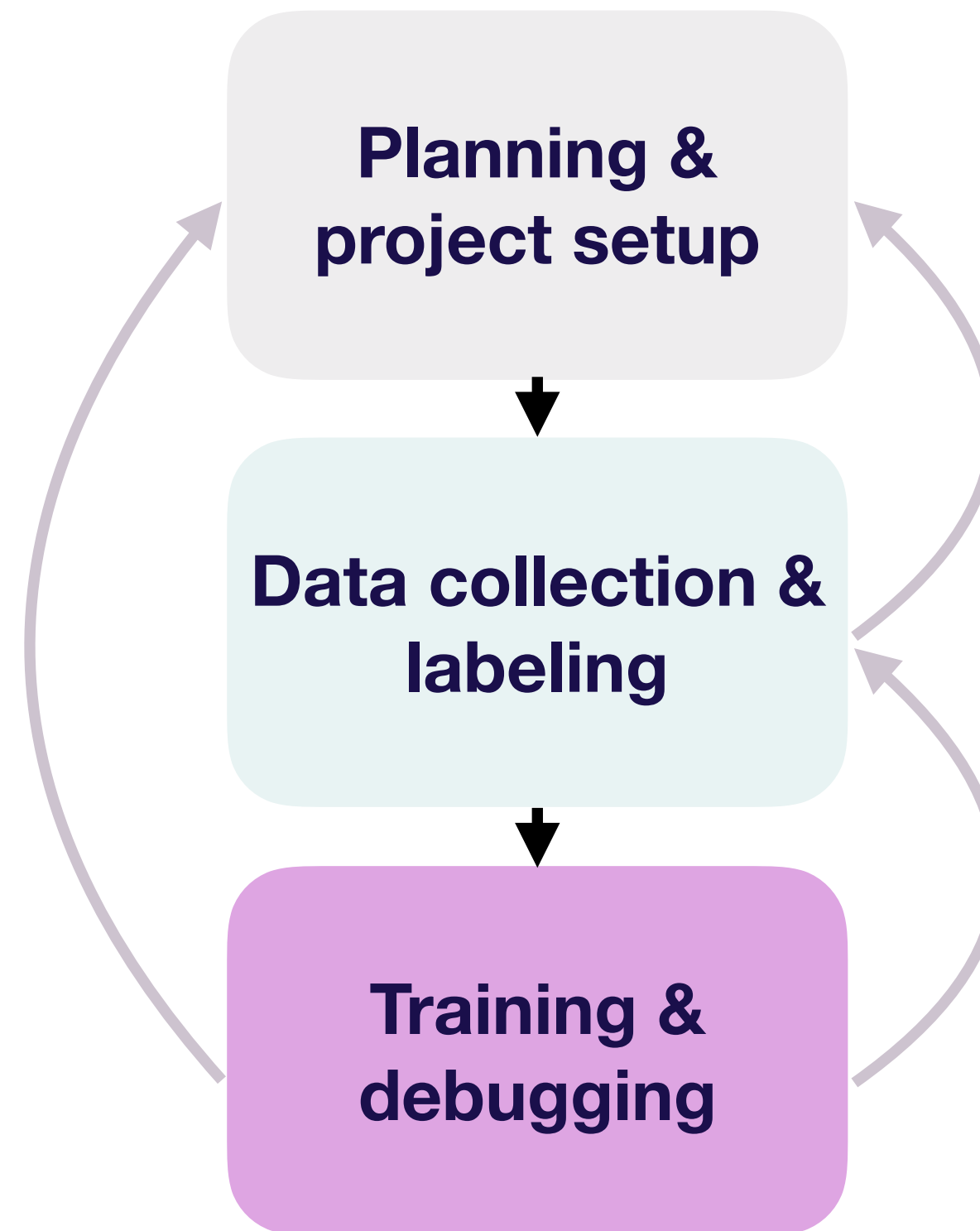
Lifecycle of a ML project



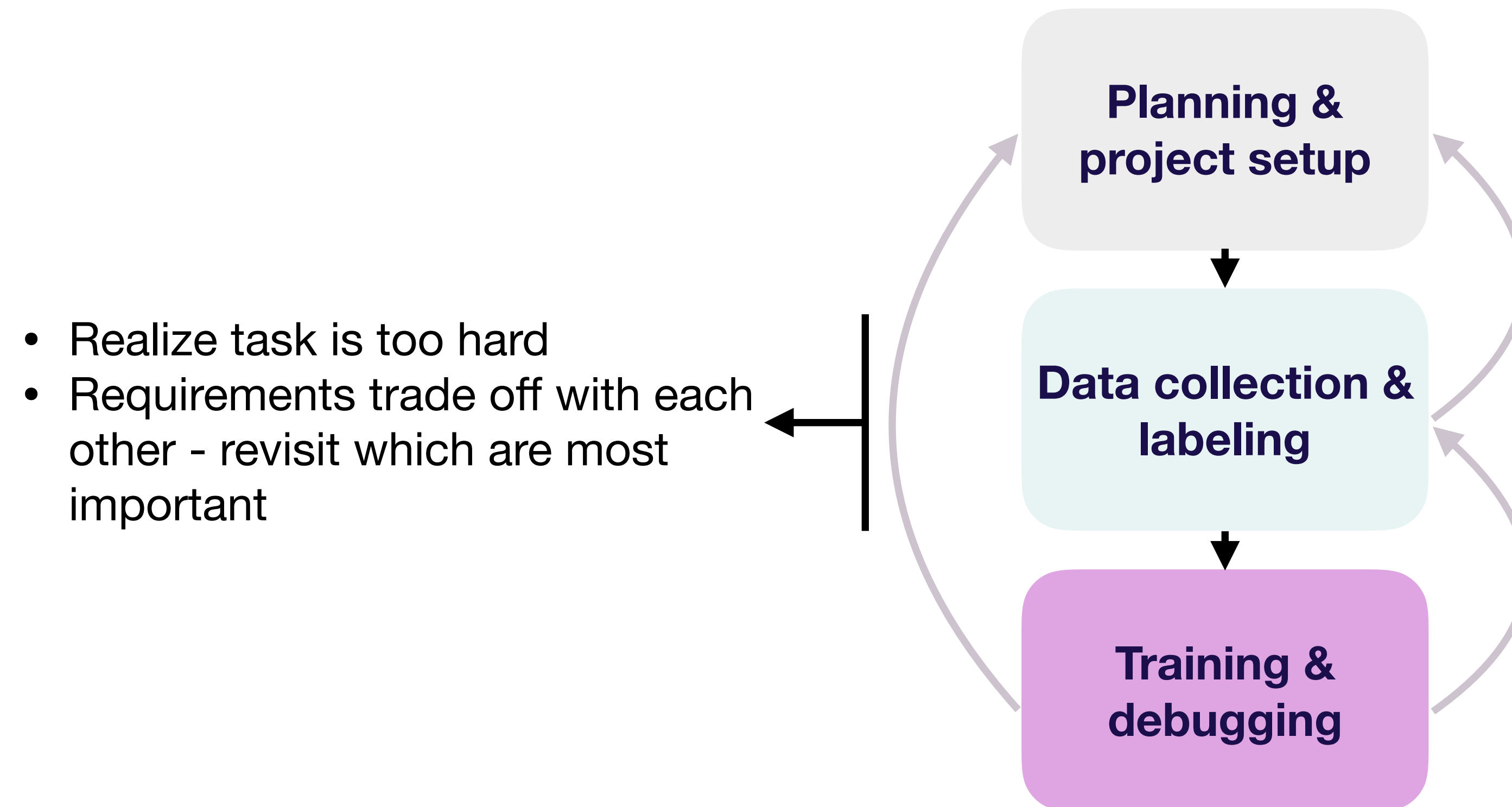
Lifecycle of a ML project



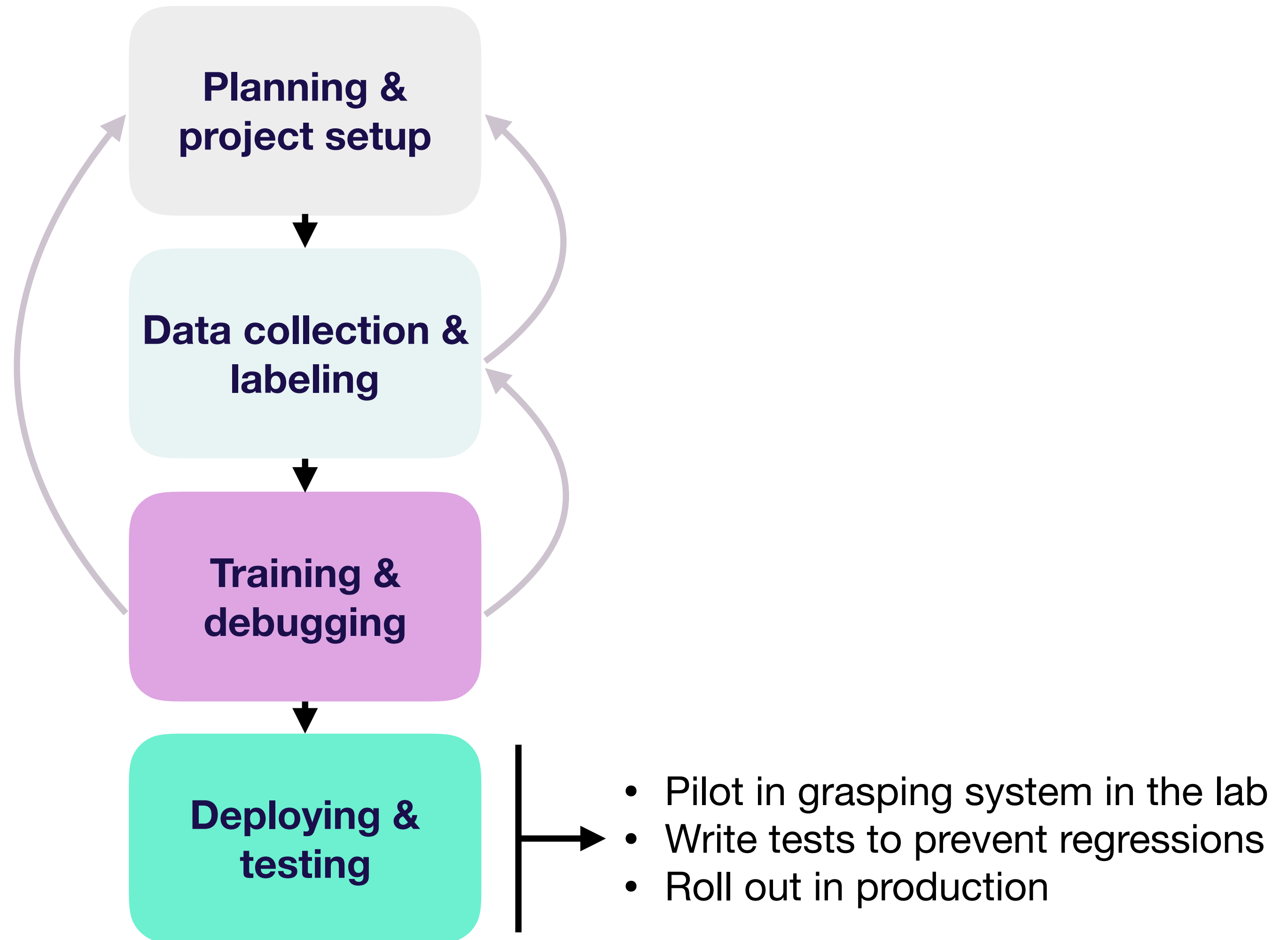
Lifecycle of a ML project



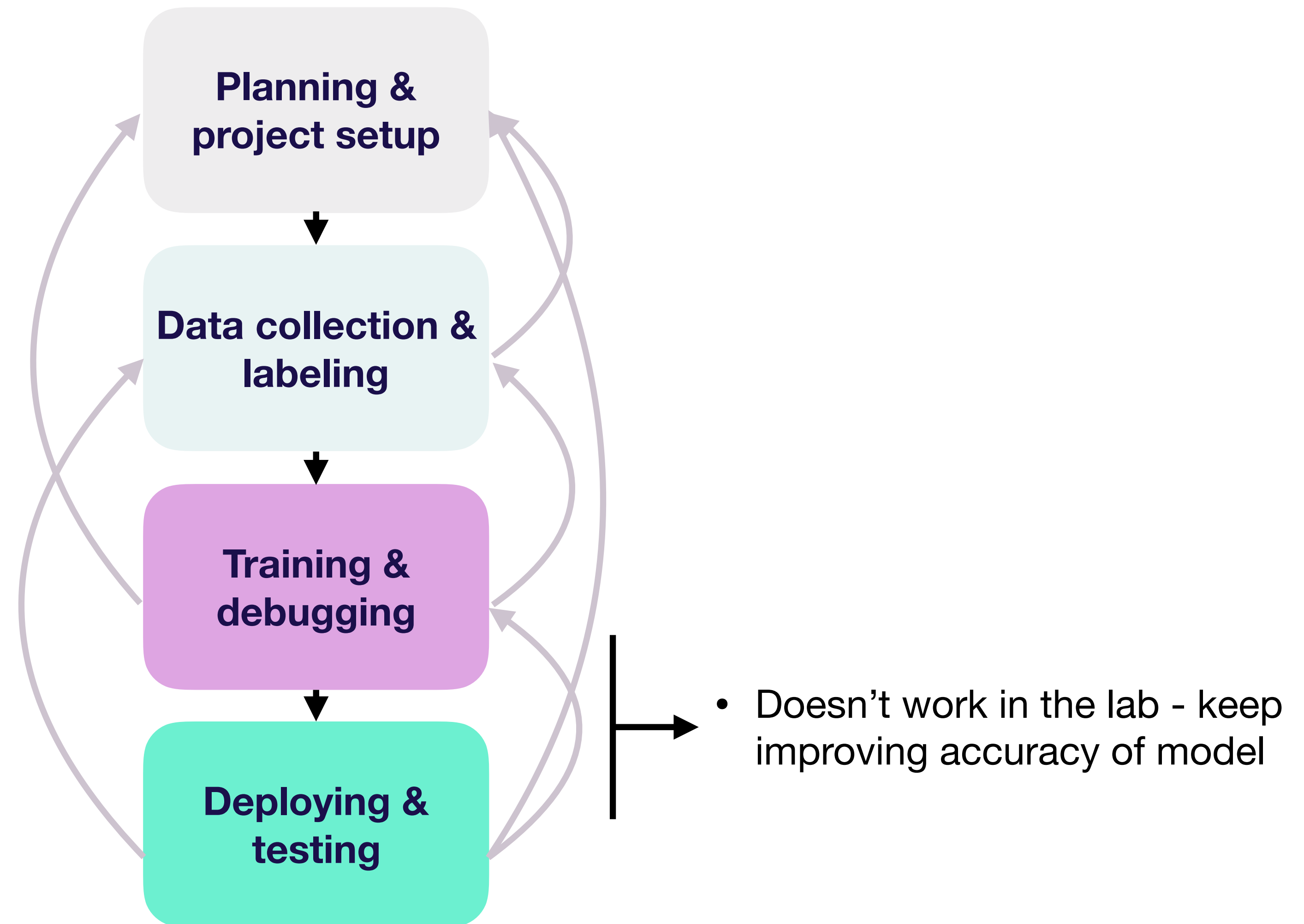
Lifecycle of a ML project



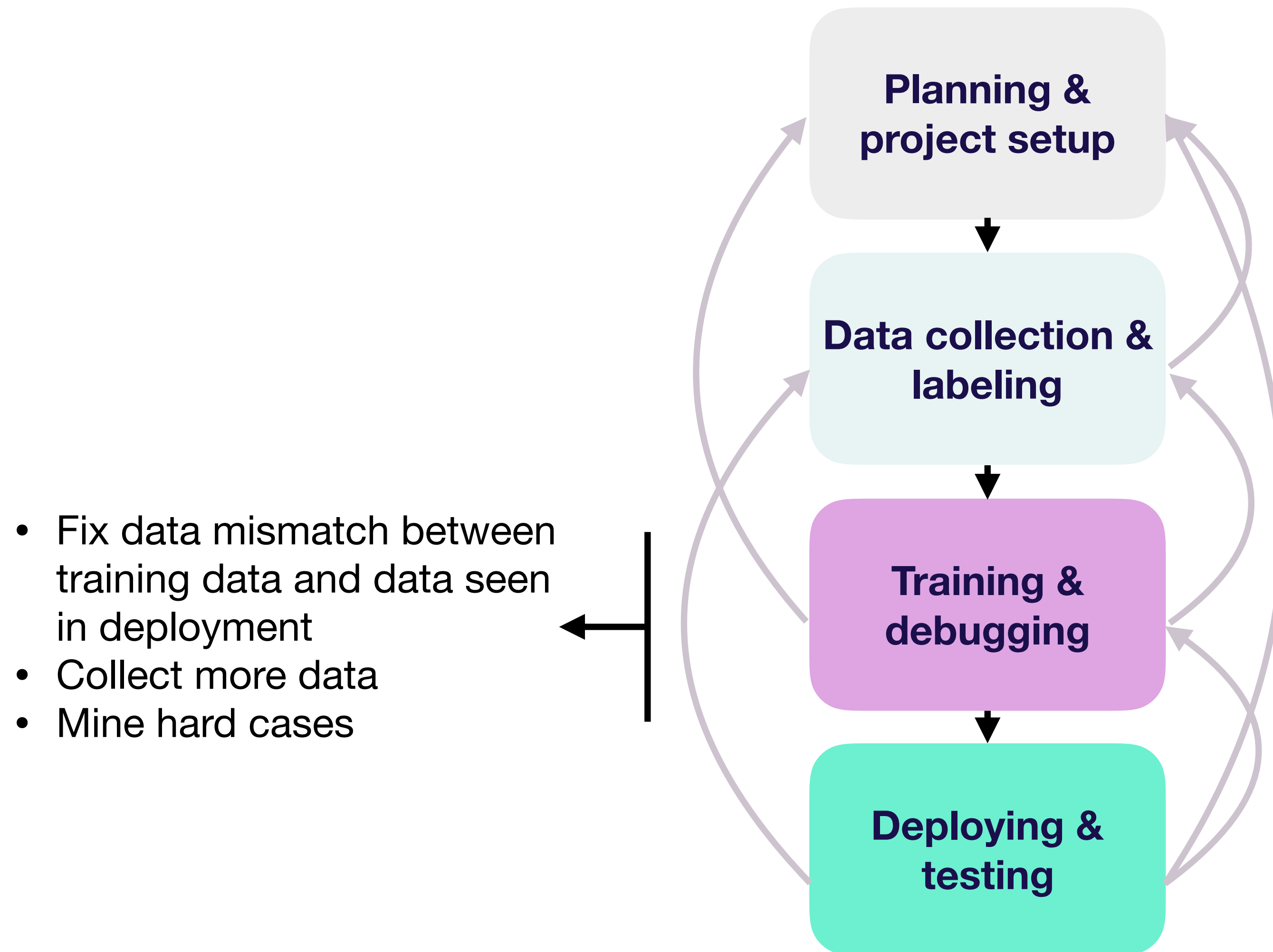
Lifecycle of a ML project



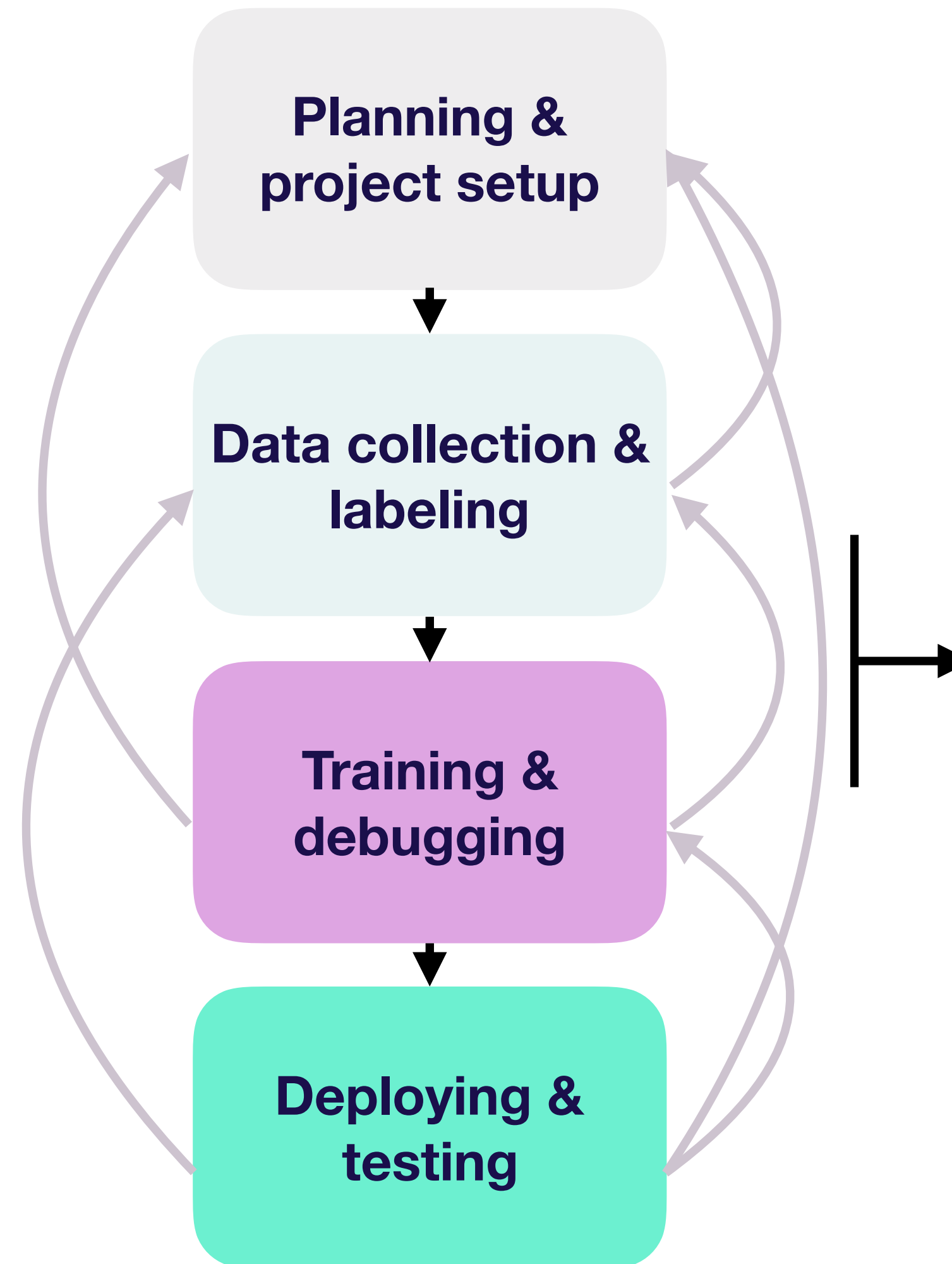
Lifecycle of a ML project



Lifecycle of a ML project

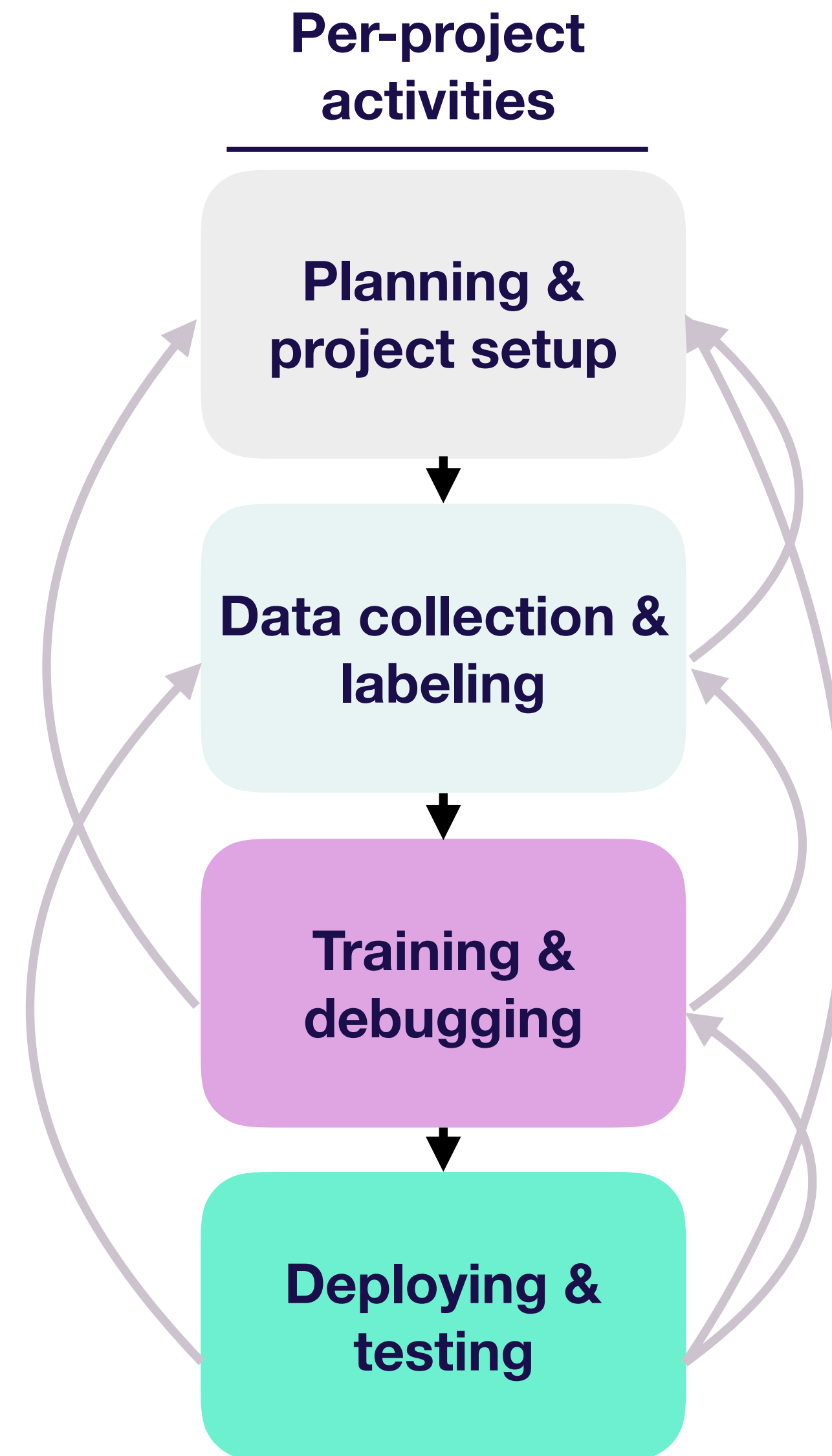


Lifecycle of a ML project

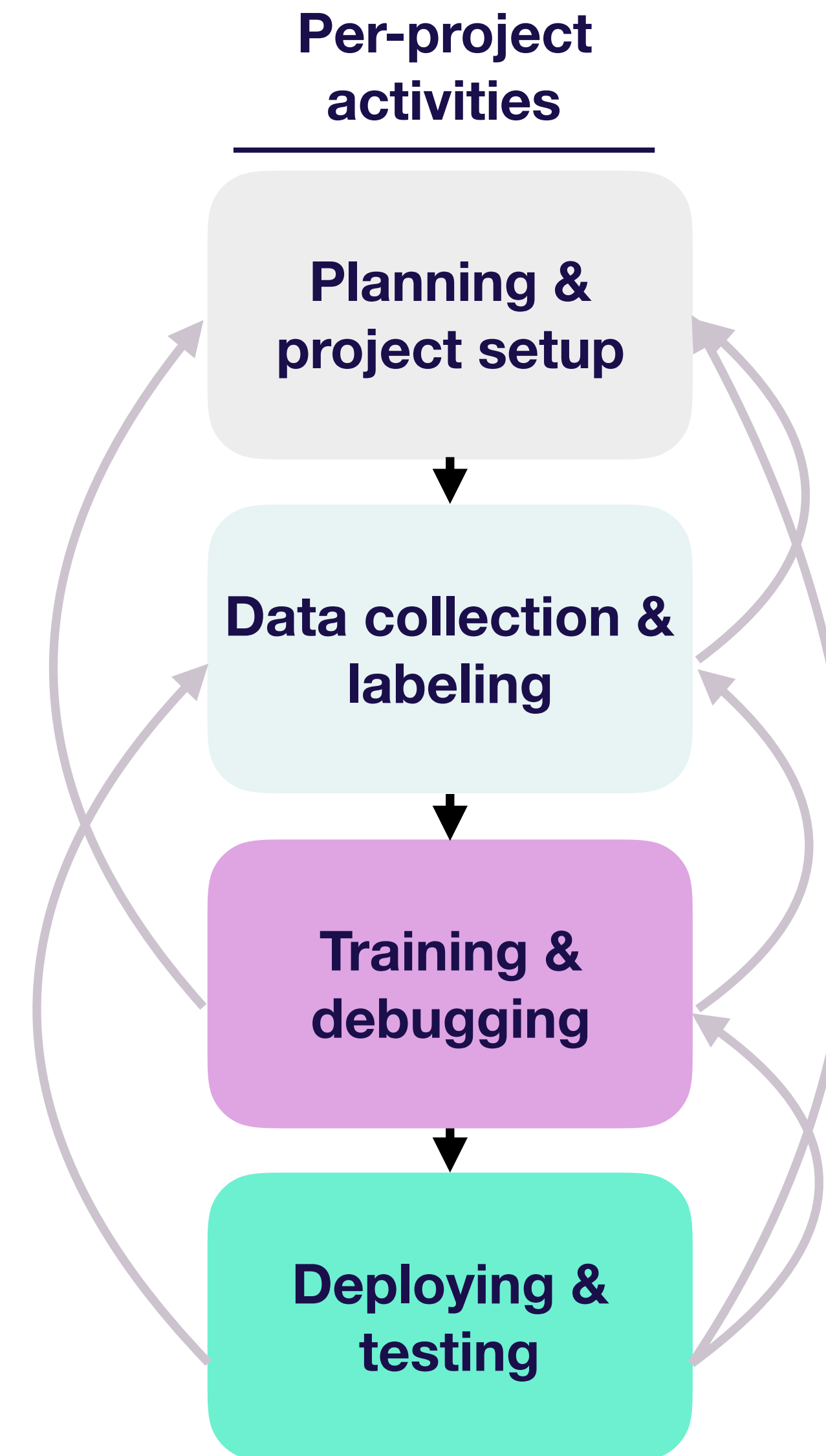


- The metric you picked doesn't actually drive downstream user behavior. Revisit the metric.
- Performance in the real world isn't great - revisit requirements (e.g., do we need to be faster or more accurate?)

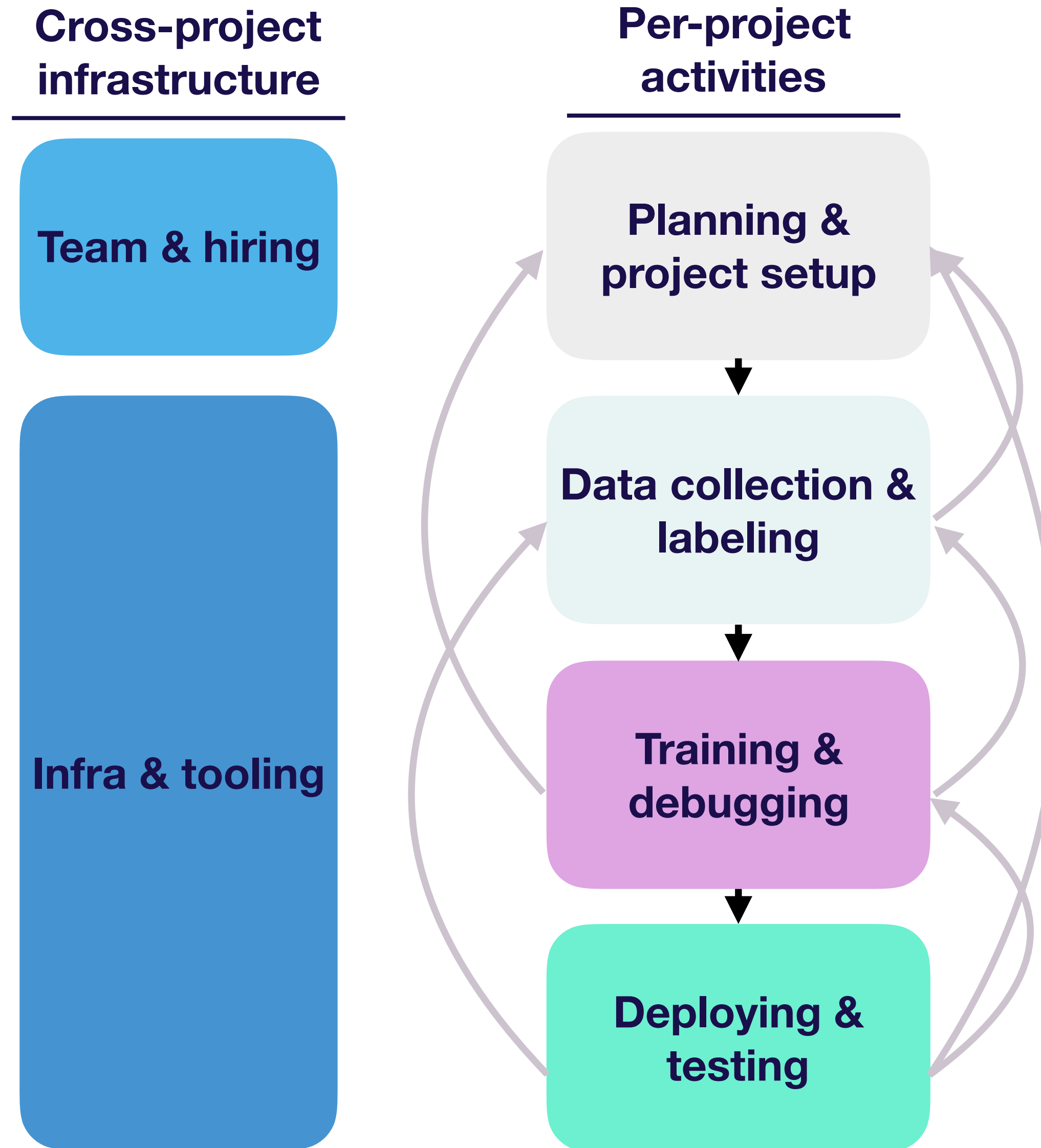
Lifecycle of a ML project



Lifecycle of a ML project



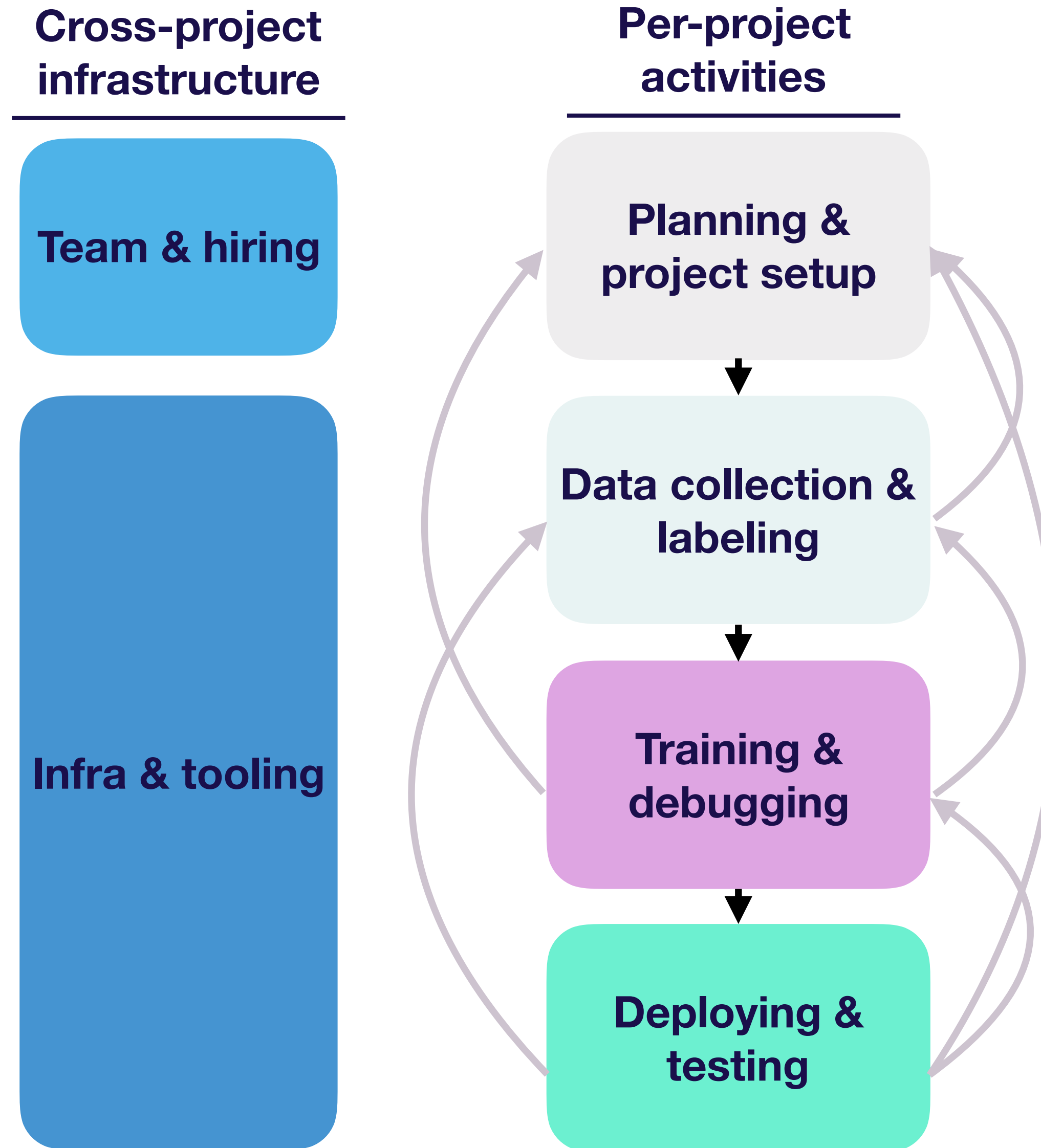
Lifecycle of a ML project



What else do you need to know?

- Understand state of the art in your domain
 - Understand what's possible
 - Know what to try next
- We will introduce most promising research areas

Lifecycle of a ML project



Questions?

Module overview

-  Lifecycle
 - How to think about all of the activities in an ML project
-  **Prioritizing projects**
 - **Assessing the feasibility and impact of your projects**
-  Archetypes
 - The main categories of ML projects, and the implications for project management
-  Metrics
 - How to pick a single number to optimize
-  Baselines
 - How to know if your model is performing well

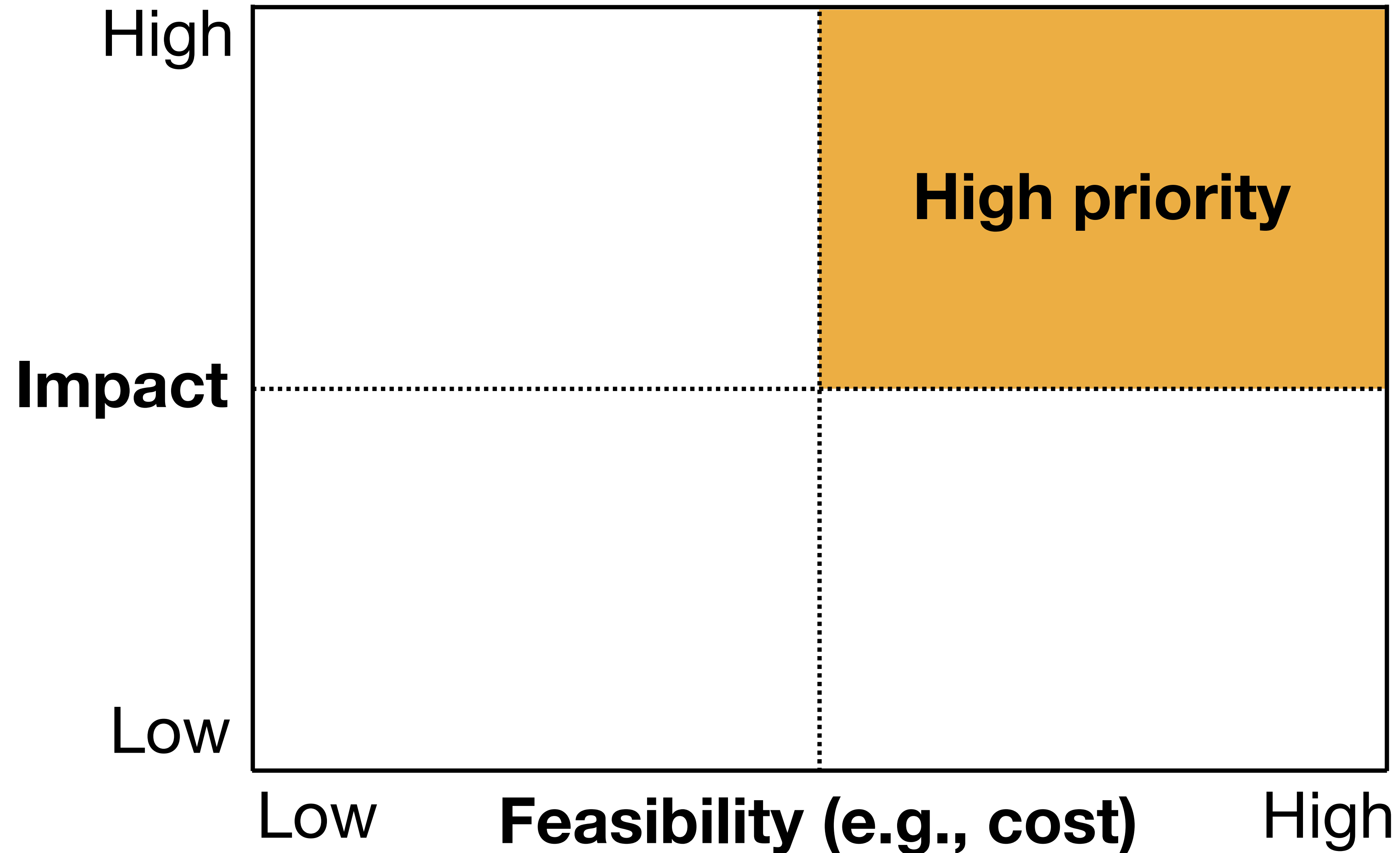
Key points for prioritizing projects

A. High-impact ML problems

- Complex parts of your pipeline
- Places where cheap prediction is valuable

B. Cost of ML projects is driven by data availability, but accuracy requirement also plays a big role

A (general) framework for prioritizing projects



Mental models for high-impact ML projects

1. Where can you take advantage of cheap prediction?
2. Where can you automate complicated manual processes?

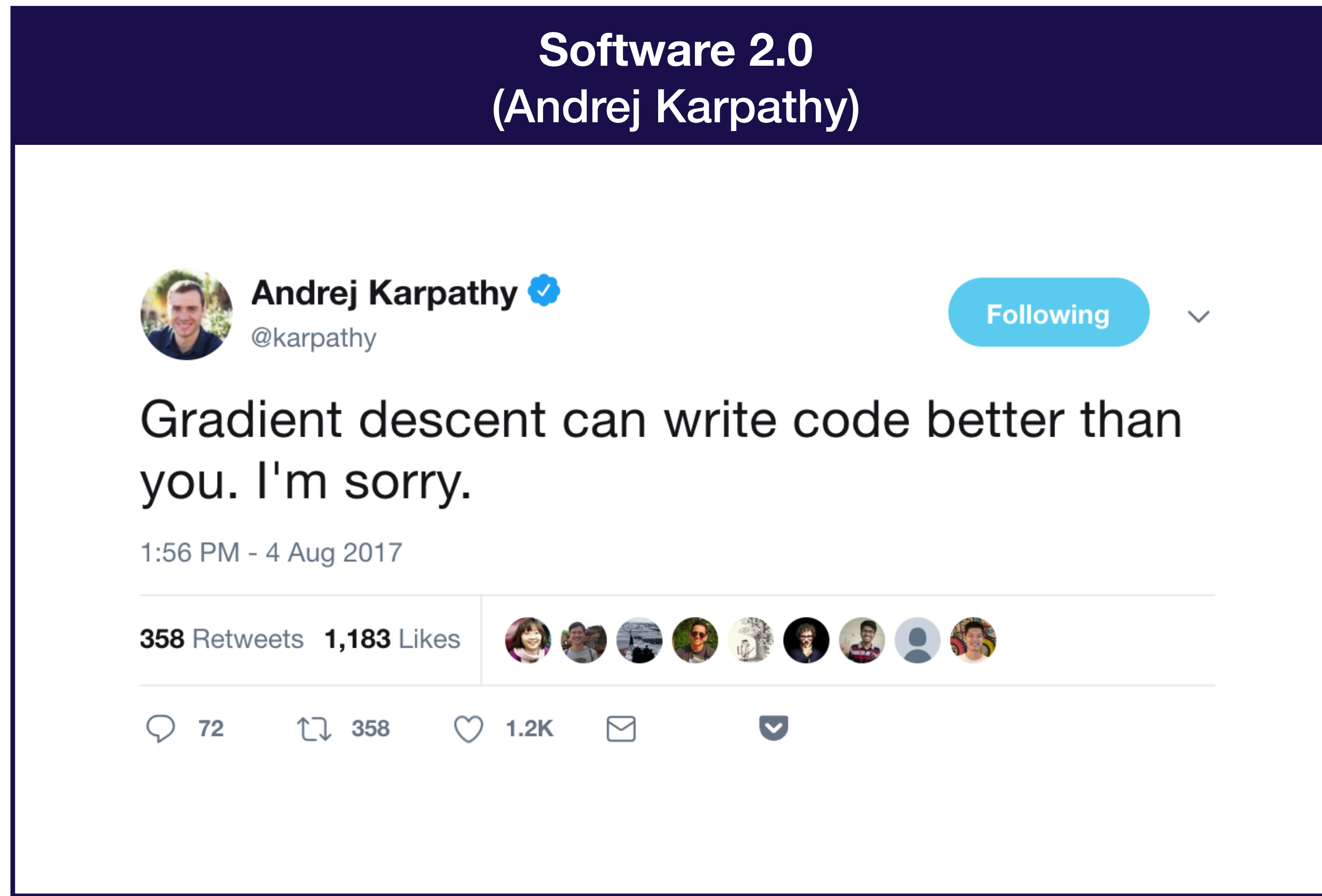
Mental models for high-impact ML projects

The economics of AI (Agrawal, Gans, Goldfarb)

- AI reduces cost of prediction
- Prediction is central for decision making
- Cheap prediction means
 - Prediction will be everywhere
 - Even in problems where it was too expensive before (e.g., for most people, hiring a driver)
- **Implication:** Look for projects where cheap prediction will have a huge business impact

Prediction Machines: The Simple Economics of Artificial Intelligence (Agrawal, Gans, Goldfarb)

Mental models for high-impact ML projects



Software 2.0 (Andrej Karpathy): <https://medium.com/@karpathy/software-2-0-a64152b37c35>

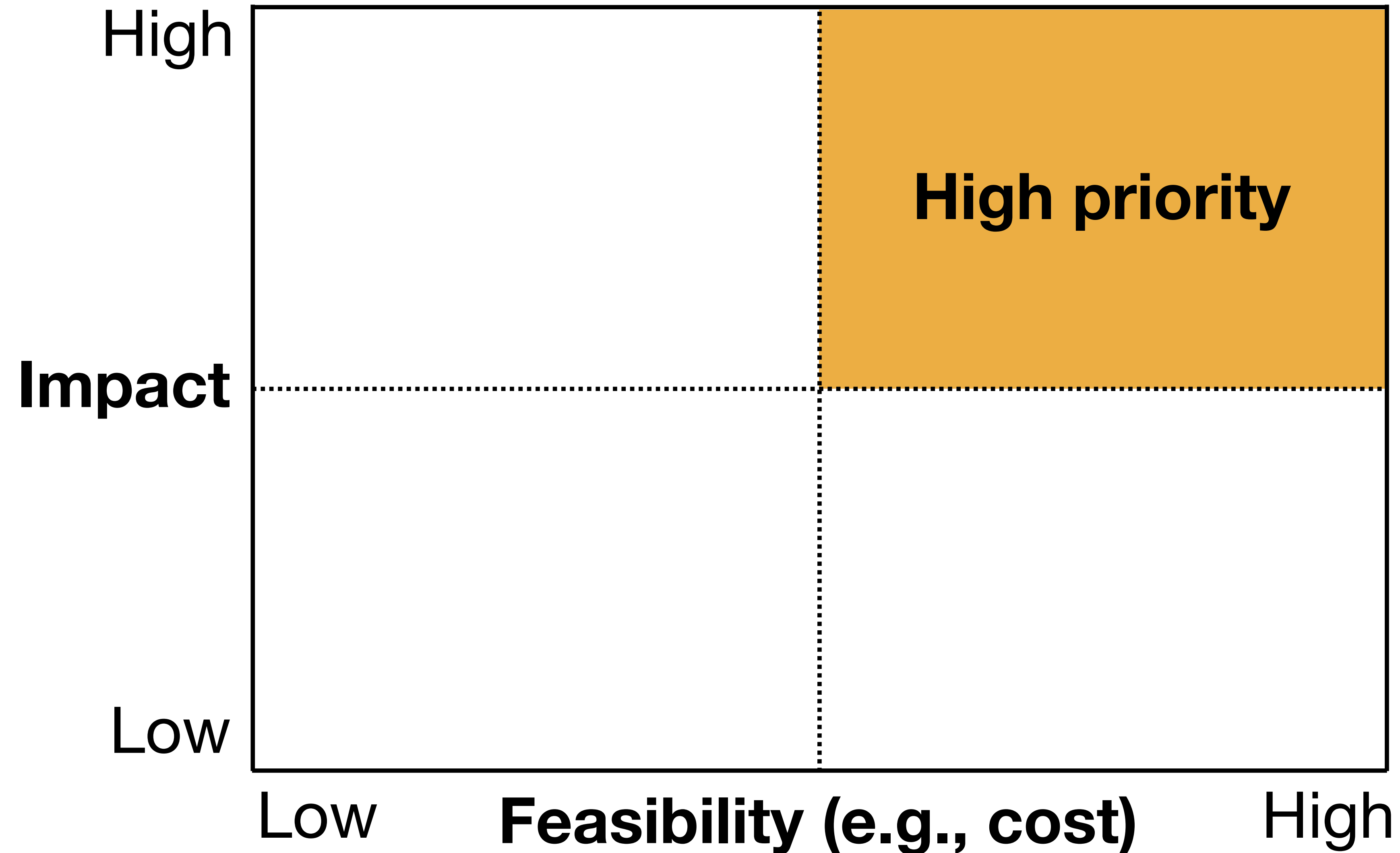
Mental models for high-impact ML projects

Software 2.0 (Andrej Karpathy)

- *Software 1.0* = traditional programs with explicit instructions (python / c++ / etc)
- Software 2.0 = humans specify goals, and algorithm searches for a program that works
- 2.0 programmers work with datasets, which get compiled via optimization
- Why? Works better, more general, computational advantages
- **Implication:** look for complicated rule-based software where we can learn the rules instead of programming them

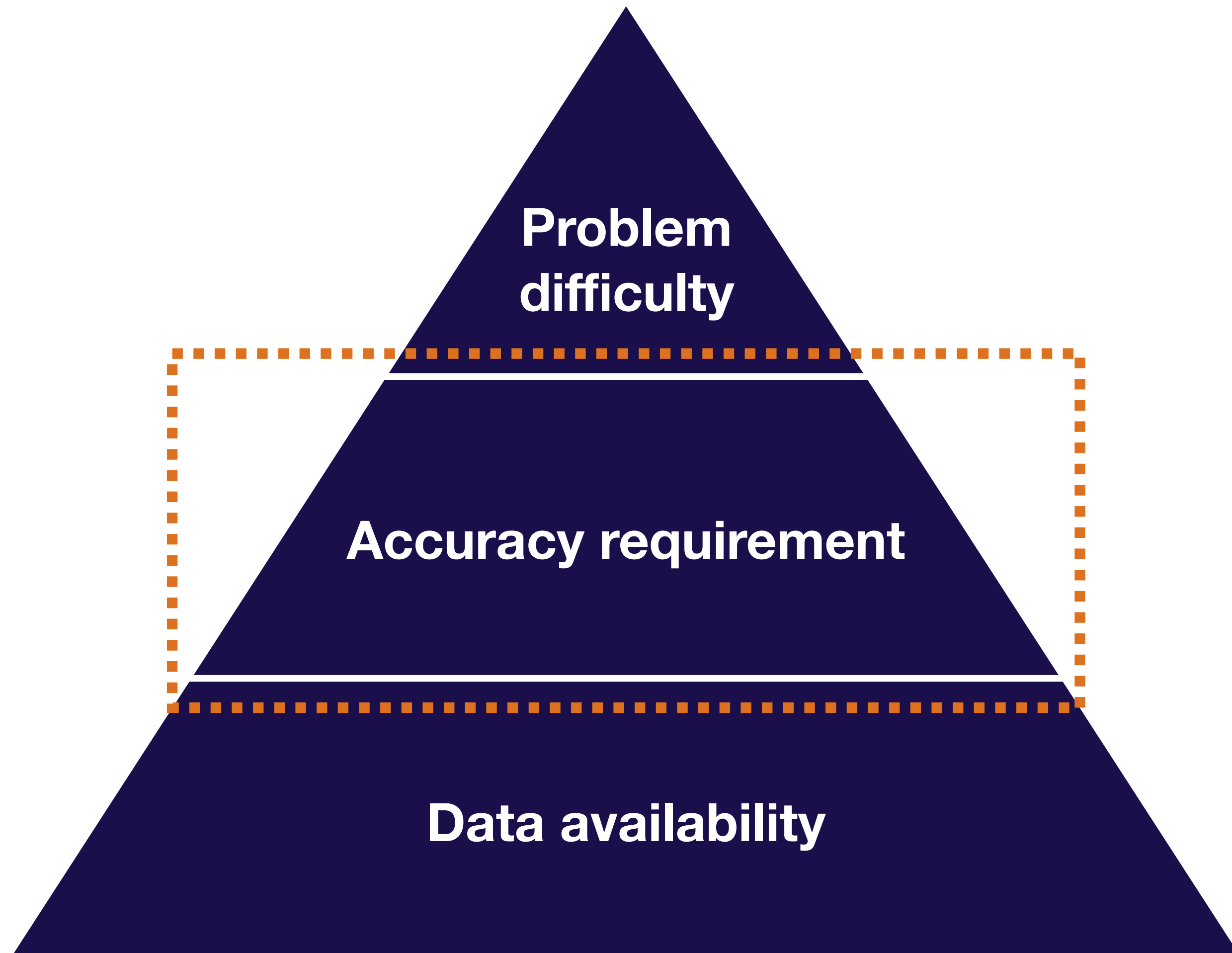
Software 2.0 (Andrej Karpathy): <https://medium.com/@karpathy/software-2-0-a64152b37c35>

A (general) framework for prioritizing projects



Assessing feasibility of ML projects

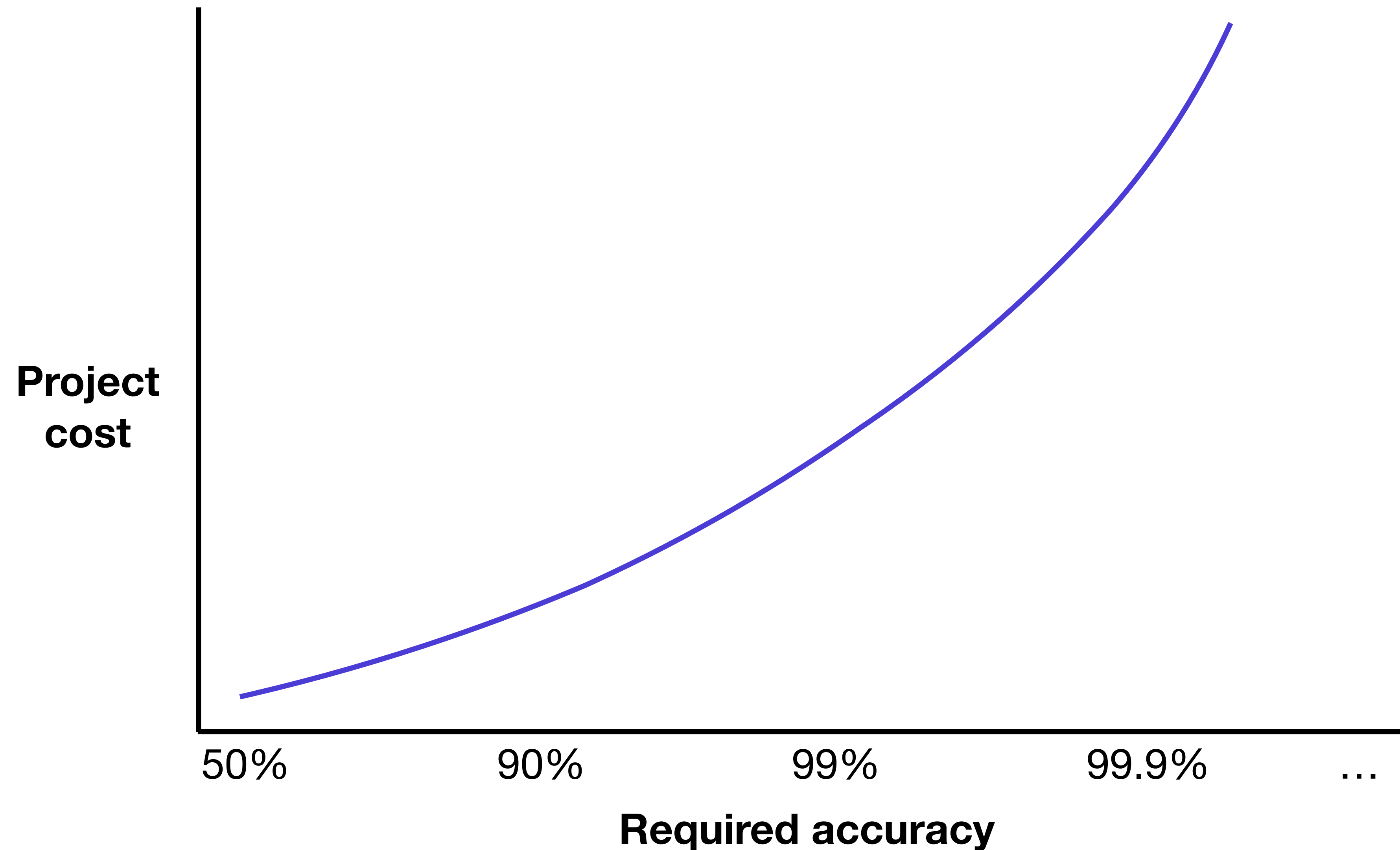
Cost drivers



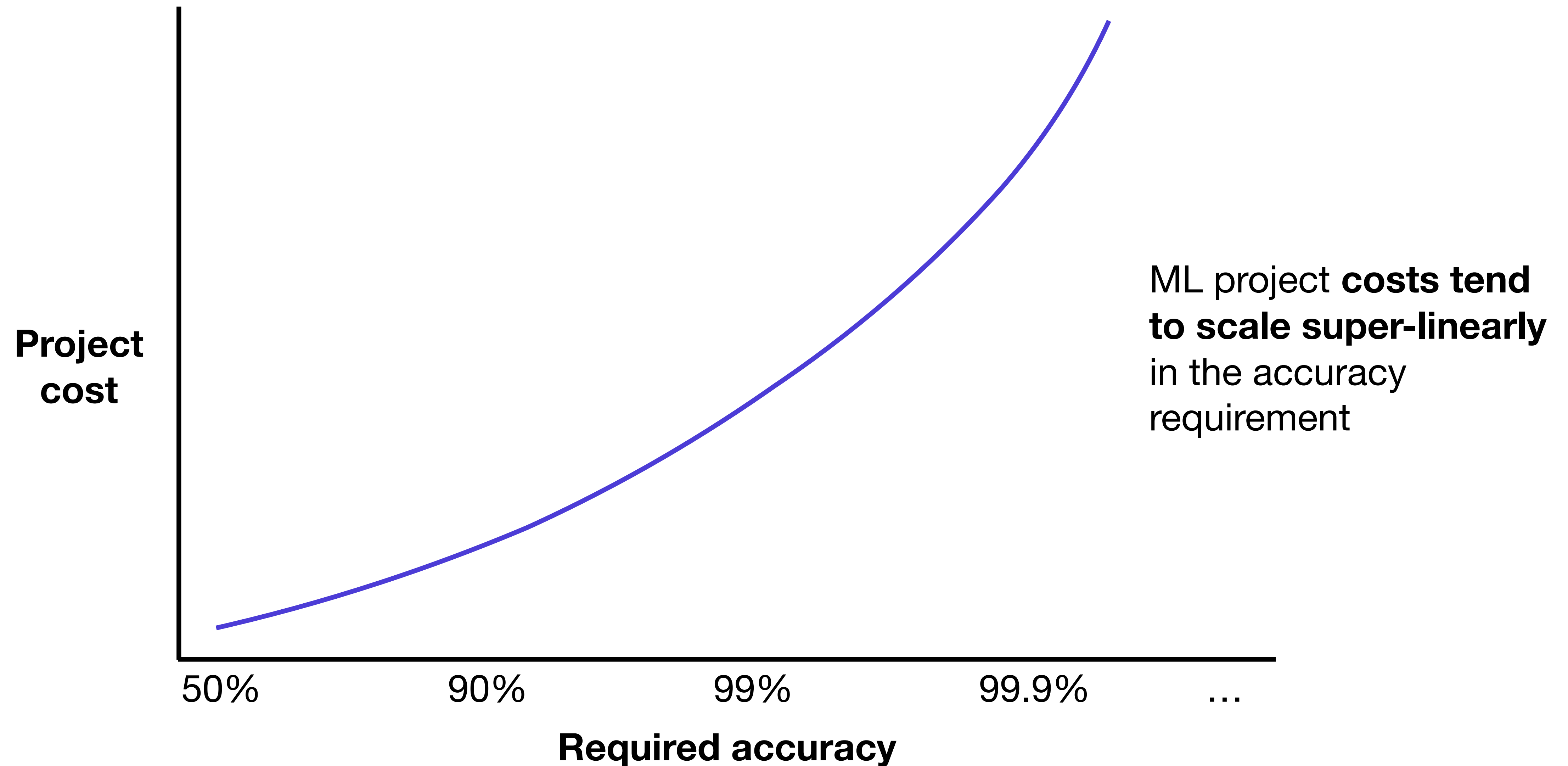
Main considerations

- Good published work on similar problems? (newer problems mean more risk & more technical effort)
 - Compute needed for training?
 - Compute available for deployment?
-
- How costly are wrong predictions?
 - How frequently does the system need to be right to be useful?
-
- How hard is it to acquire data?
 - How expensive is data labeling?
 - How much data will be needed?

Why are accuracy requirements so important?

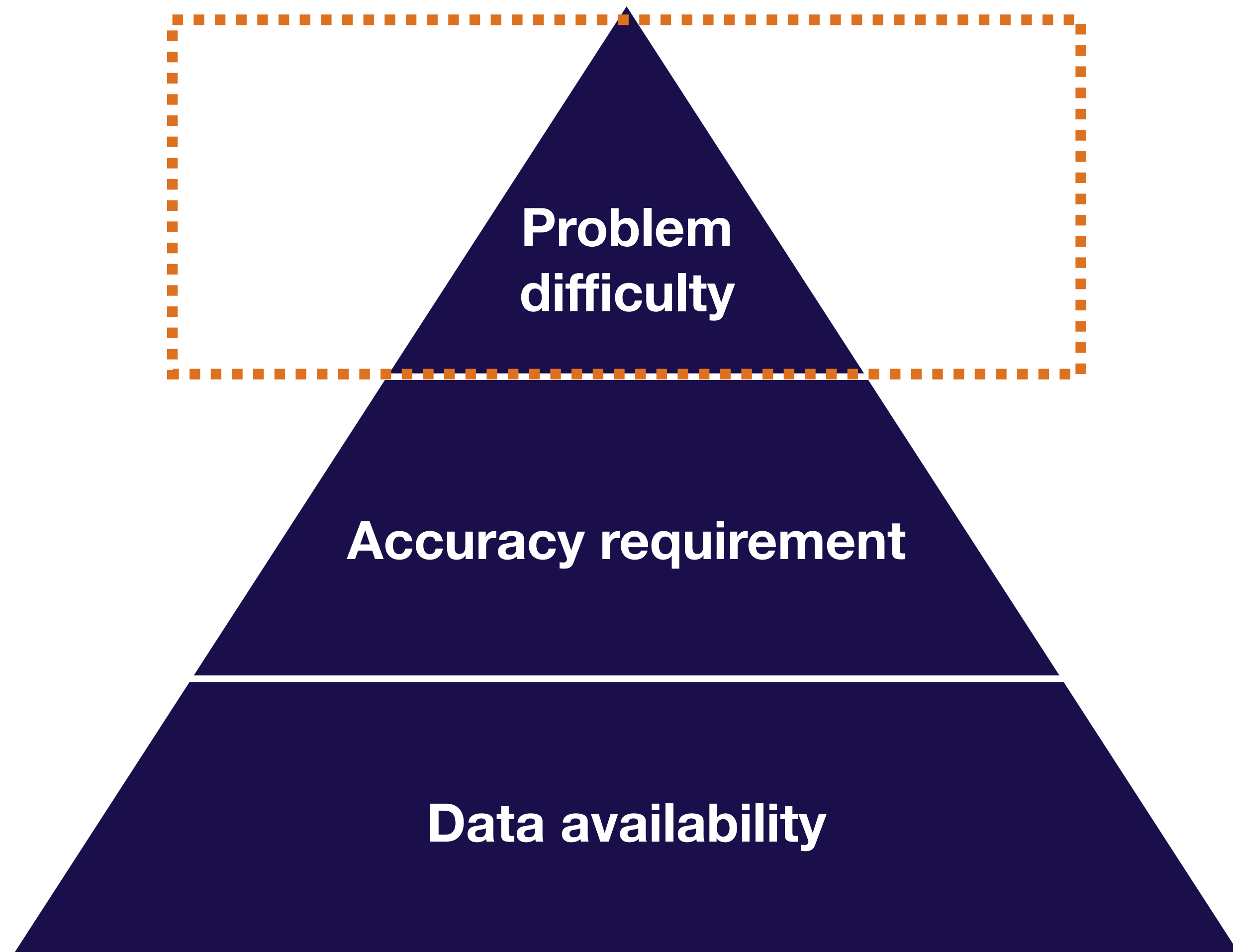


Why are accuracy requirements so important?



Assessing feasibility of ML projects

Cost drivers



Main considerations

- Good published work on similar problems? (newer problems mean more risk & more technical effort)
 - Compute needed for training?
 - Compute available for deployment?
-
- How costly are wrong predictions?
 - How frequently does the system need to be right to be useful?
-
- How hard is it to acquire data?
 - How expensive is data labeling?
 - How much data will be needed?

What's still hard in machine learning?

"It may be a hundred years before a computer beats humans at Go -- maybe even longer," said Dr. Piet Hut, an astrophysicist at the Institute for Advanced Study in Princeton, N.J., and a fan of the game. "If a reasonably intelligent person learned to play Go, in a few months he could beat all existing computer programs. You don't have to be a Kasparov."

New York Times, July 1997

What's still hard in machine learning?



What's still hard in machine learning?



Following

Pretty much anything that a normal person can do in <1 sec, we can now automate with AI.

Examples

- Recognize content of images
- Understand speech
- Translate speech
- Grasp objects
- etc.

Counter-examples?

- Understand humor / sarcasm
- In-hand robotic manipulation
- Generalize to new scenarios
- etc.

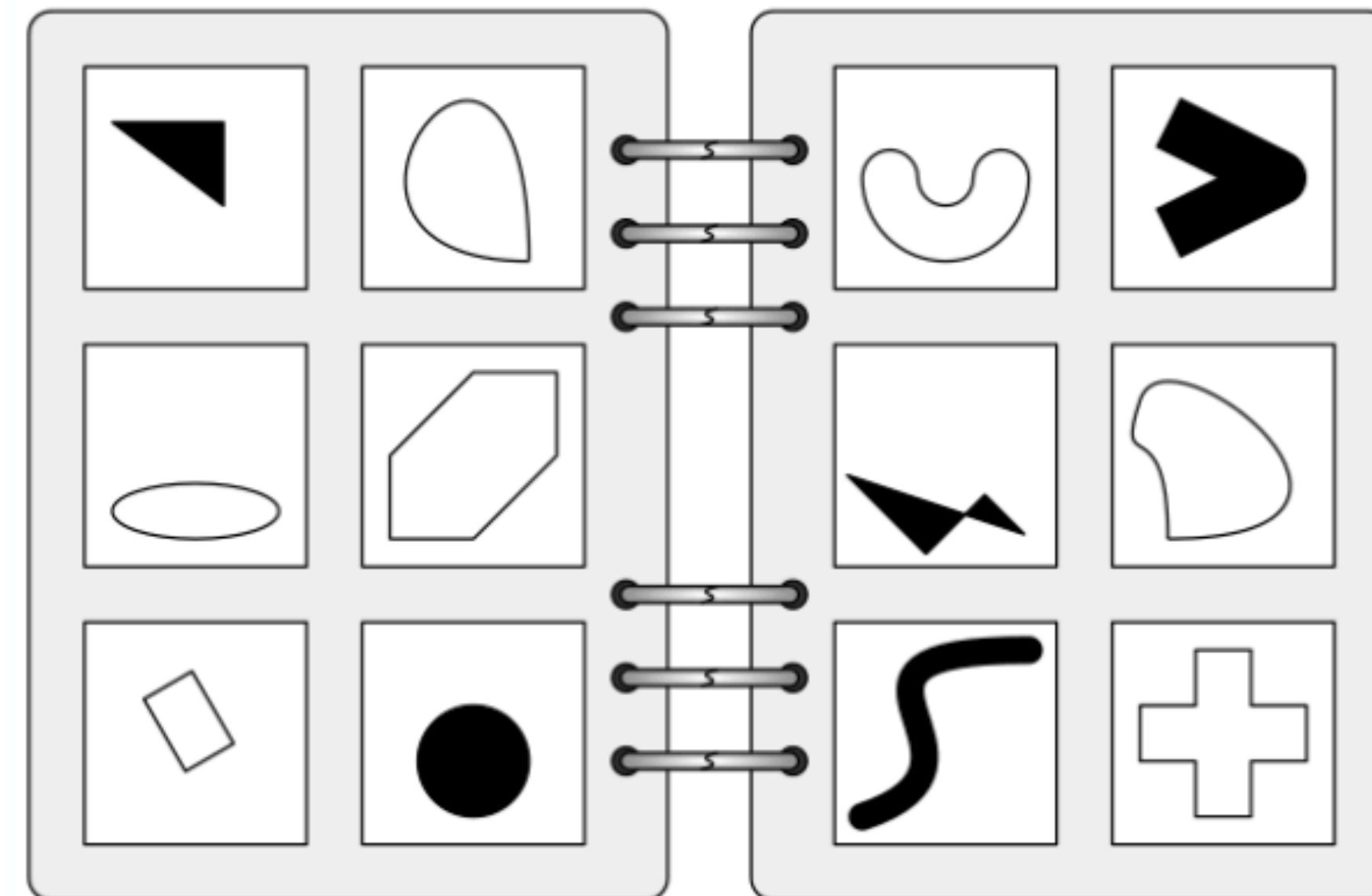
What's still hard in machine learning?

- Unsupervised learning
- Reinforcement learning
- Both are showing promise in limited domains where tons of data and compute are available

What's still hard in supervised learning?

- Answering questions
- Summarizing text
- Predicting video
- Building 3D models
- Real-world speech recognition
- Resisting adversarial examples
- Doing math
- Solving word puzzles
- Bongard problems
- Etc

Example of a Bongard Problem



What types of problems are hard?

	<u>Instances</u>	<u>Examples</u>
Output is complex	<ul style="list-style-type: none">● High-dimensional output● Ambiguous output	<ul style="list-style-type: none">● 3D reconstruction● Video prediction● Dialog systems
Reliability is required	<ul style="list-style-type: none">● High precision is required● Robustness is required	<ul style="list-style-type: none">● Failing safely out-of-distribution● Robustness to adversarial attacks● High-precision pose estimation
Generalization is required	<ul style="list-style-type: none">● Out of distribution data● Reasoning, planning, causality	<ul style="list-style-type: none">● Self-driving: edge cases● Self-driving: control● Small data



Why is FSR focusing on pose estimation?

Impact

- FSR's goal is grasping - requires reliable pose estimation
- Traditional robotics pipeline uses hand-designed heuristics & online optimization
 - Slow
 - Brittle
 - Great candidate for Software 2.0!

Feasibility






- Data availability
 - Easy to collect data
 - Labeling data could be a challenge, but can instrument lab with sensors
- Accuracy requirement
 - Require high accuracy to grasp an object: $<0.5\text{cm}$
 - However, low cost of failure - picks per hour important, not % successes
- Problem difficulty
 - Similar published results exist but need to adapt to our objects and robot

Key points for prioritizing projects

- A. To find high-impact ML problems, look for complex parts of your pipeline and places where cheap prediction is valuable
- B. The cost of ML projects is primarily driven by data availability, but your accuracy requirement also plays a big role

Questions?

Module overview

-  Lifecycle
 - How to think about all of the activities in an ML project
-  Prioritizing projects
 - Assessing the feasibility and impact of your projects
-  **Archetypes**
 - **The main categories of ML projects, and the implications for project management**
-  Metrics
 - How to pick a single number to optimize
-  Baselines
 - How to know if your model is performing well



Machine learning project archetypes

Improve an
existing process

Examples

- Improve code completion in an IDE
- Build a customized recommendation system
- Build a better video game AI



Machine learning project archetypes

Examples

Improve an existing process

- Improve code completion in an IDE
- Build a customized recommendation system
- Build a better video game AI

Augment a manual process

- Turn sketches into slides
- Email auto-completion
- Help a radiologist do their job faster



Machine learning project archetypes

Examples

Improve an existing process

- Improve code completion in an IDE
- Build a customized recommendation system
- Build a better video game AI

Augment a manual process

- Turn sketches into slides
- Email auto-completion
- Help a radiologist do their job faster

Automate a manual process

- Full self-driving
- Automated customer support
- Automated website design



Machine learning project archetypes

Improve an
existing process

Key questions

- Do your models truly improve performance?
- Does performance improvement generate business value?
- Do performance improvements lead to a data flywheel?



Machine learning project archetypes

Key questions

Improve an
existing process

- Do your models truly improve performance?
- Does performance improvement generate business value?
- Do performance improvements lead to a data flywheel?

Augment a manual
process

- How good does the system need to be to be useful?
- How can you collect enough data to make it that good?



Machine learning project archetypes

Key questions

Improve an existing process

- Do your models truly improve performance?
- Does performance improvement generate business value?
- Do performance improvements lead to a data flywheel?

Augment a manual process

- How good does the system need to be to be useful?
- How can you collect enough data to make it that good?

Automate a manual process

- What is an acceptable failure rate for the system?
- How can you guarantee that it won't exceed that failure rate?
- How inexpensively can you label data from the system?



Machine learning project archetypes

Key questions

Improve an existing process

- Do your models truly improve performance?
- Does performance improvement generate business value?
- Do performance improvements lead to a data flywheel?

Augment a manual process

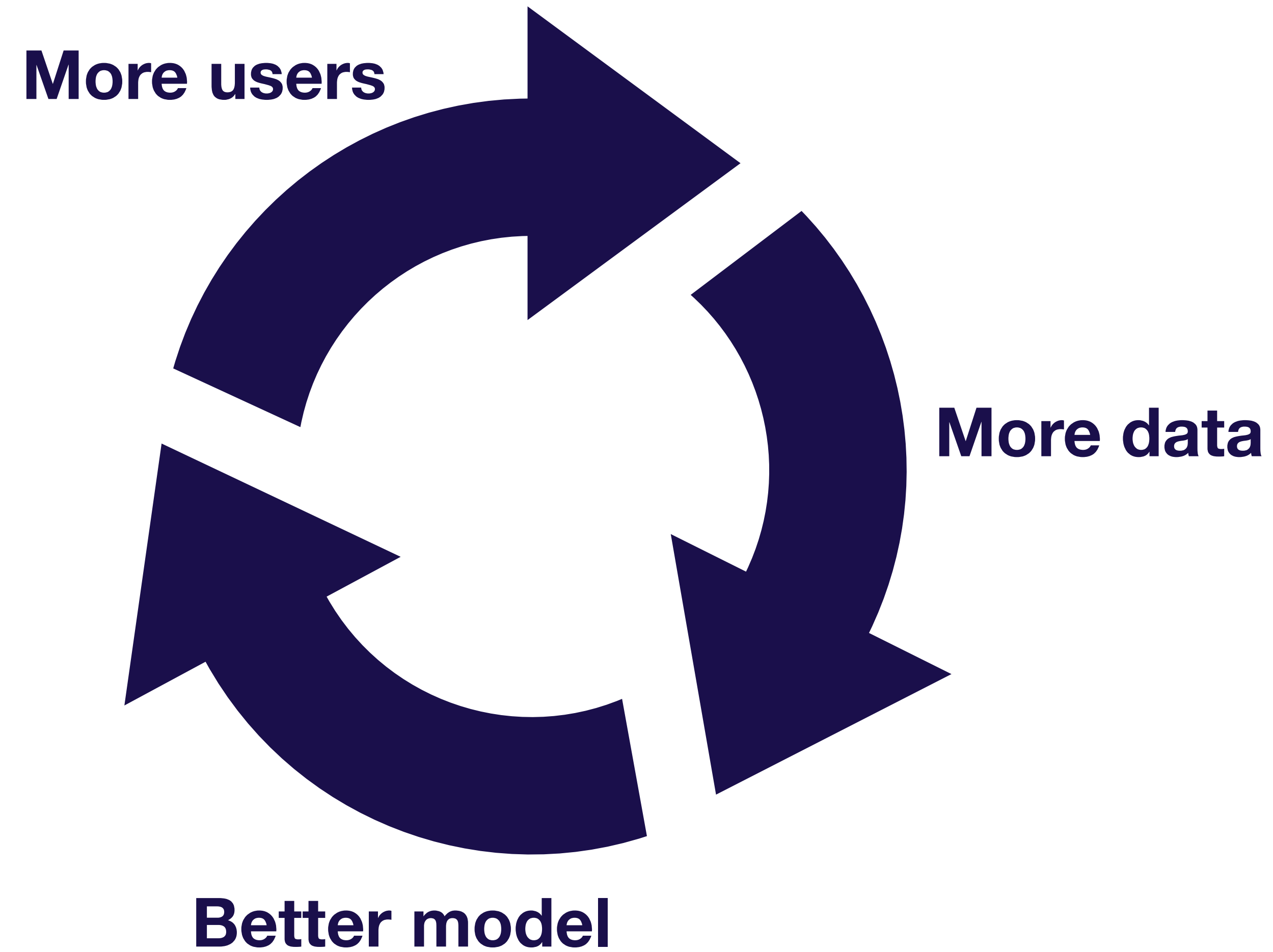
- How good does the system need to be to be useful?
- How can you collect enough data to make it that good?

Automate a manual process

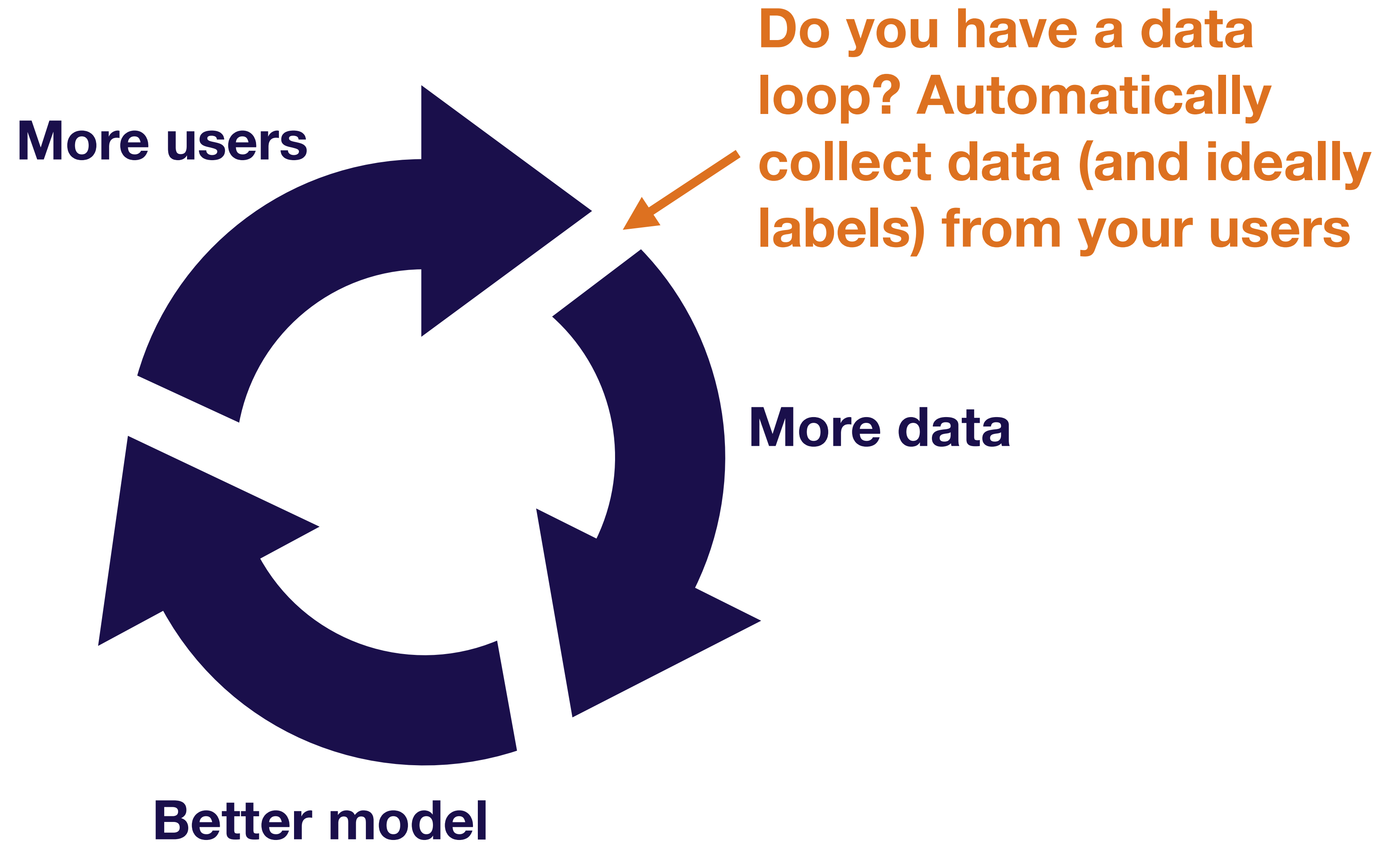
- What is an acceptable failure rate for the system?
- How can you guarantee that it won't exceed that failure rate?
- How inexpensively can you label data from the system?



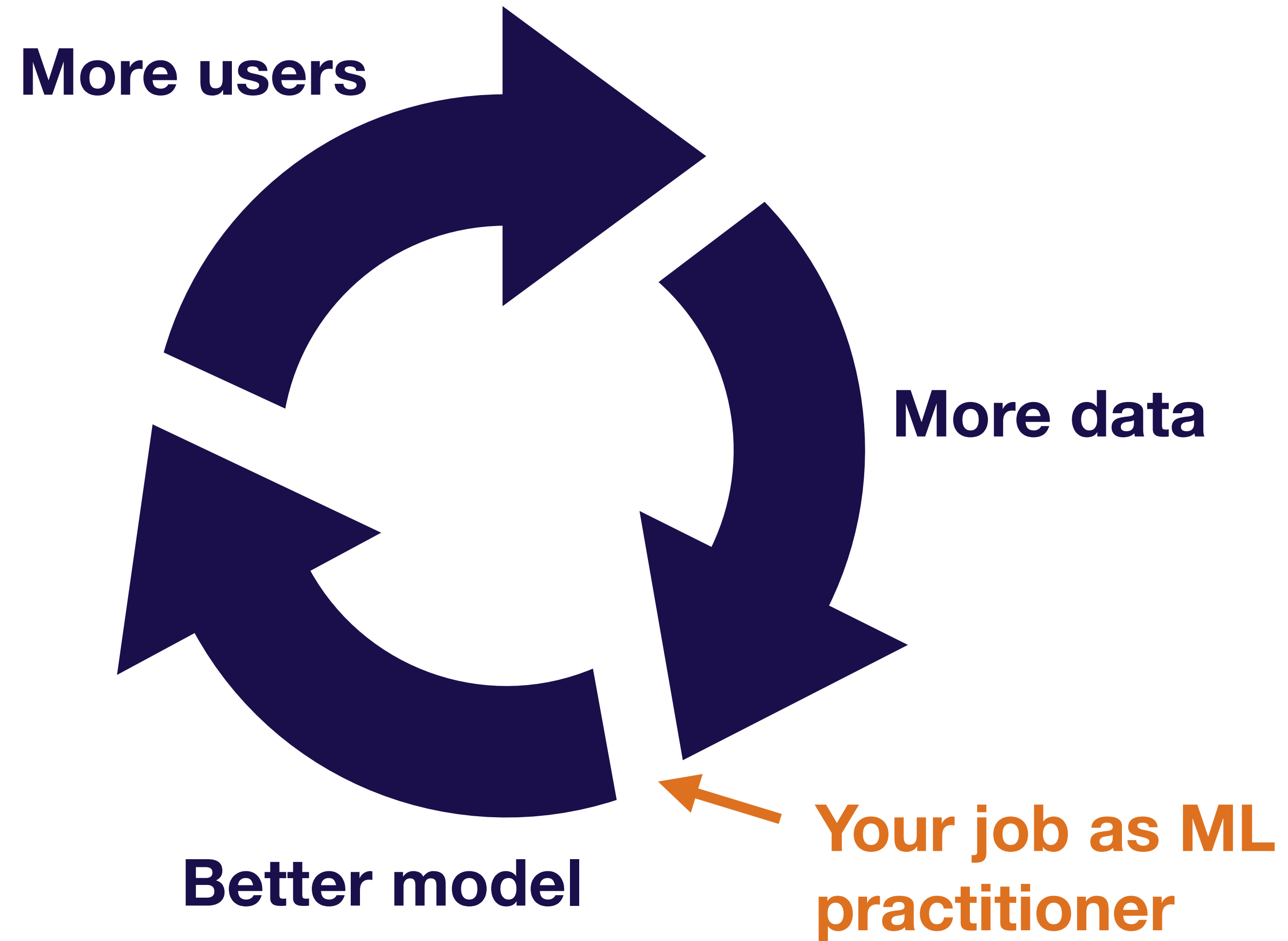
Data flywheels



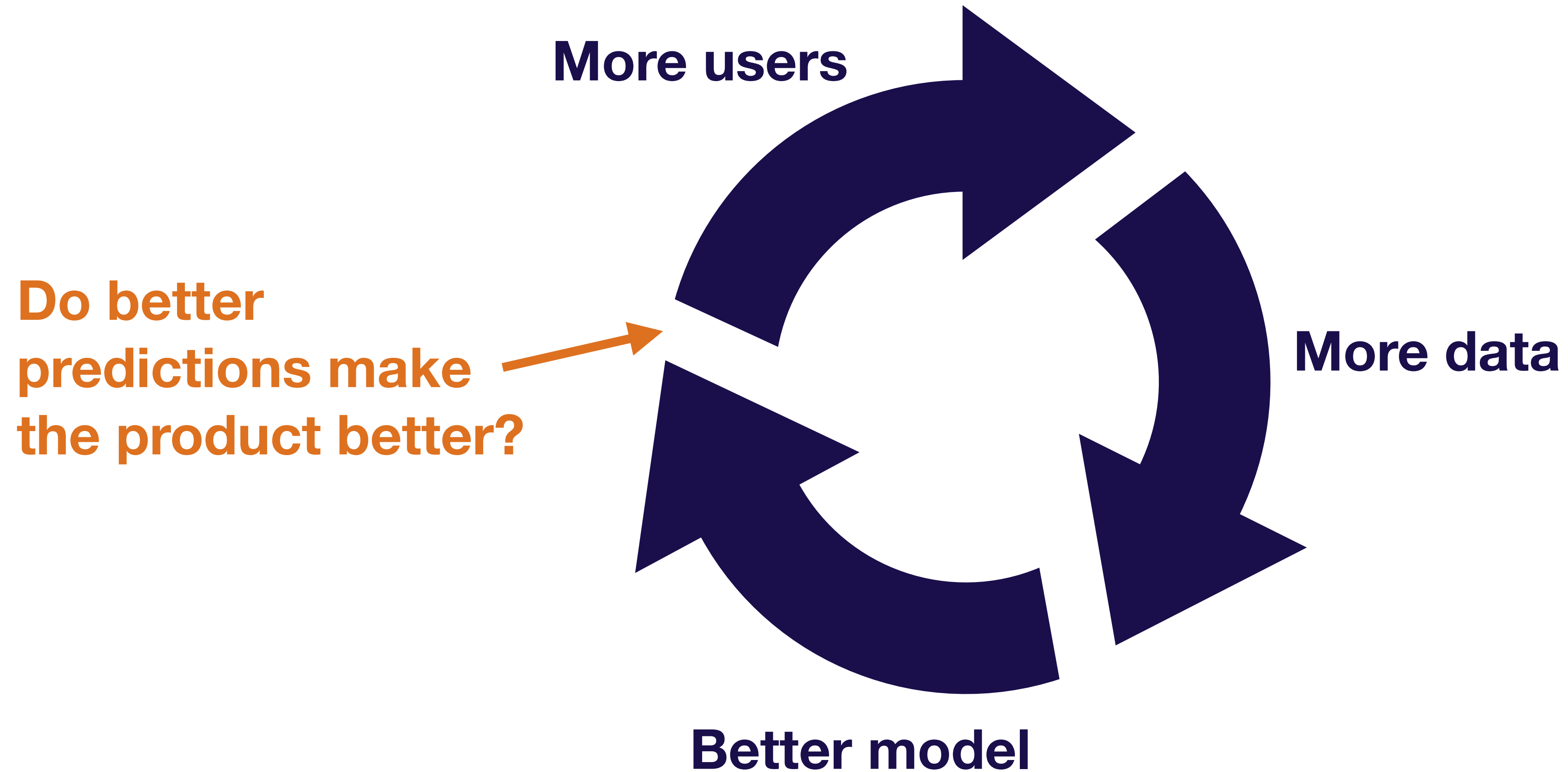
Data flywheels



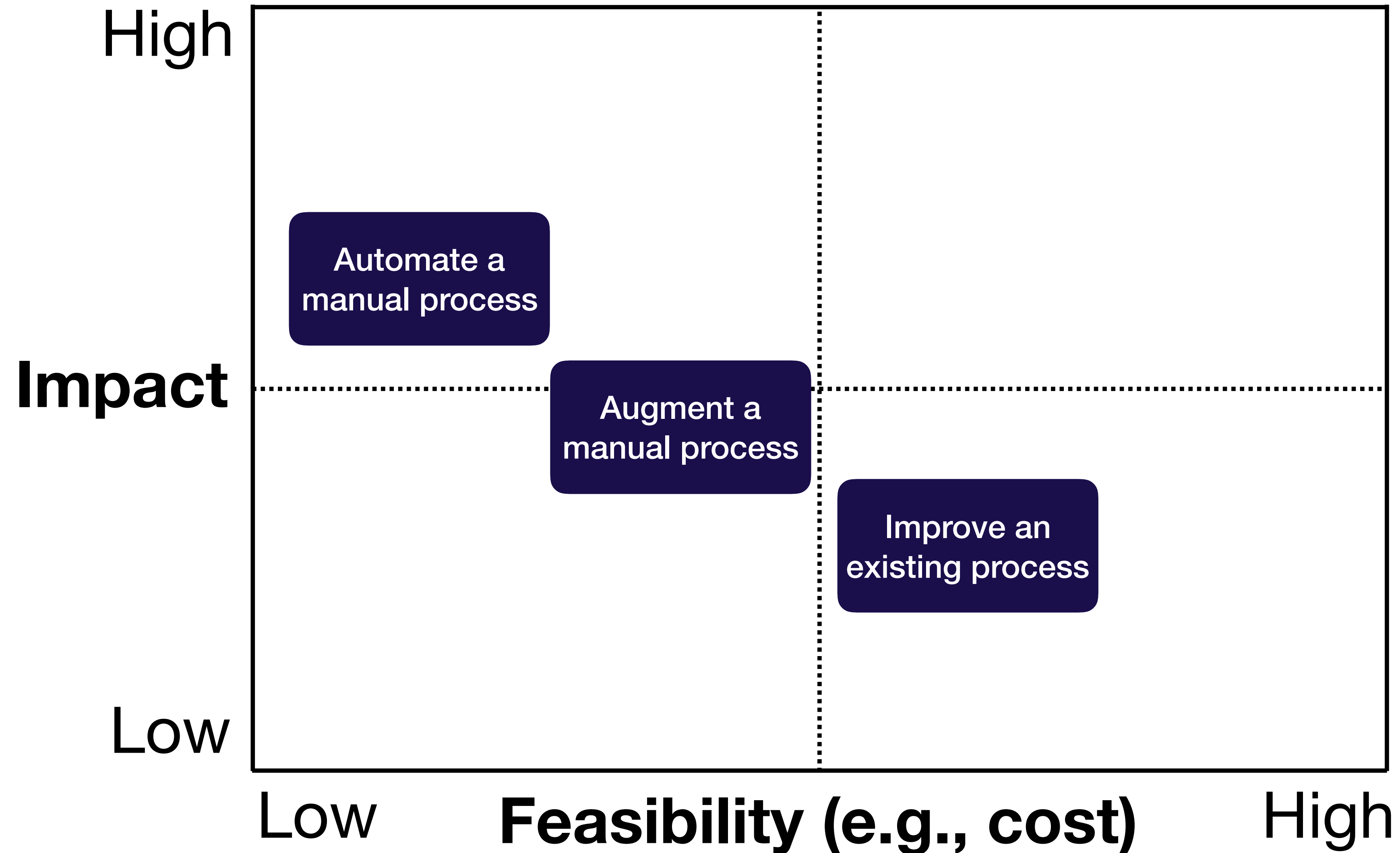
Data flywheels



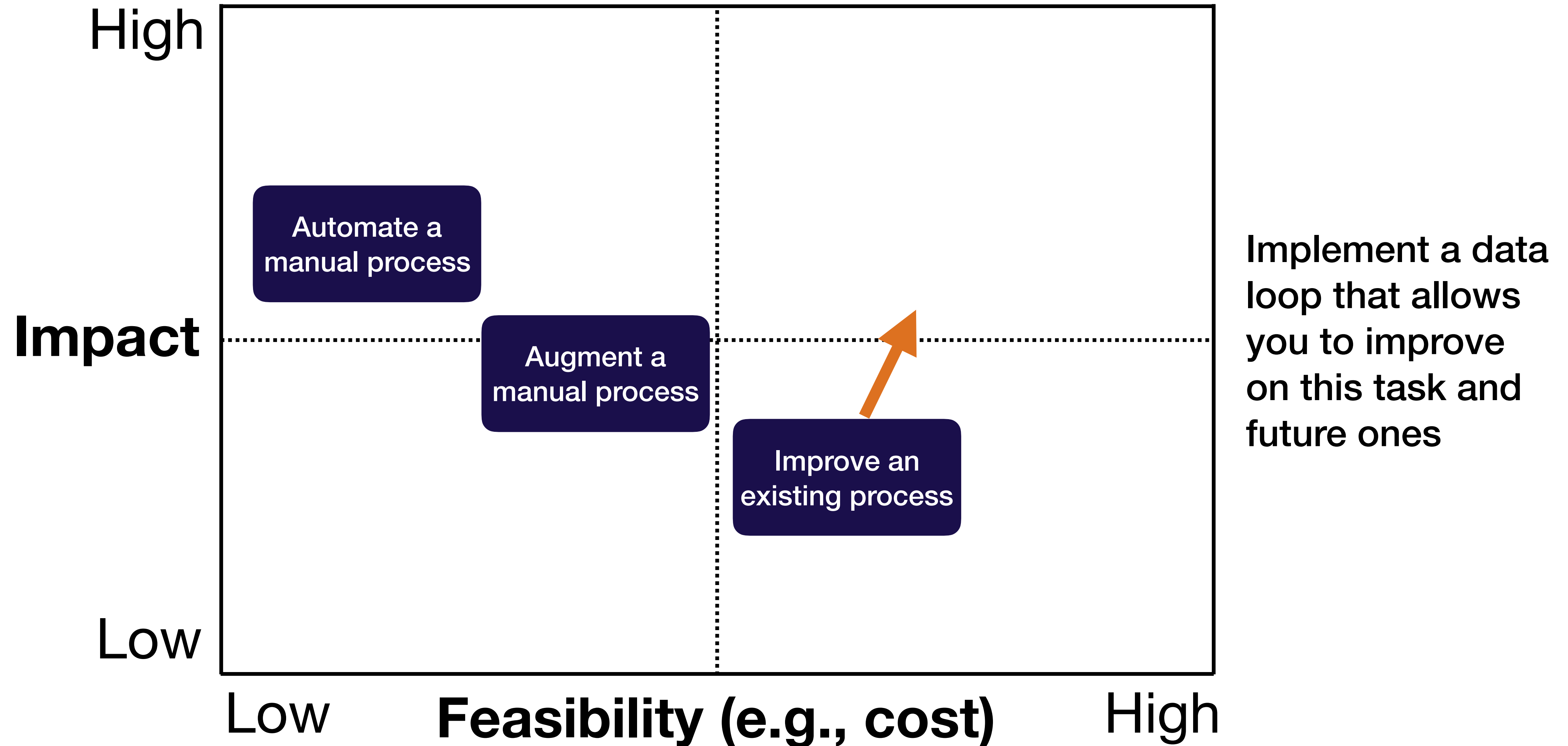
Data flywheels



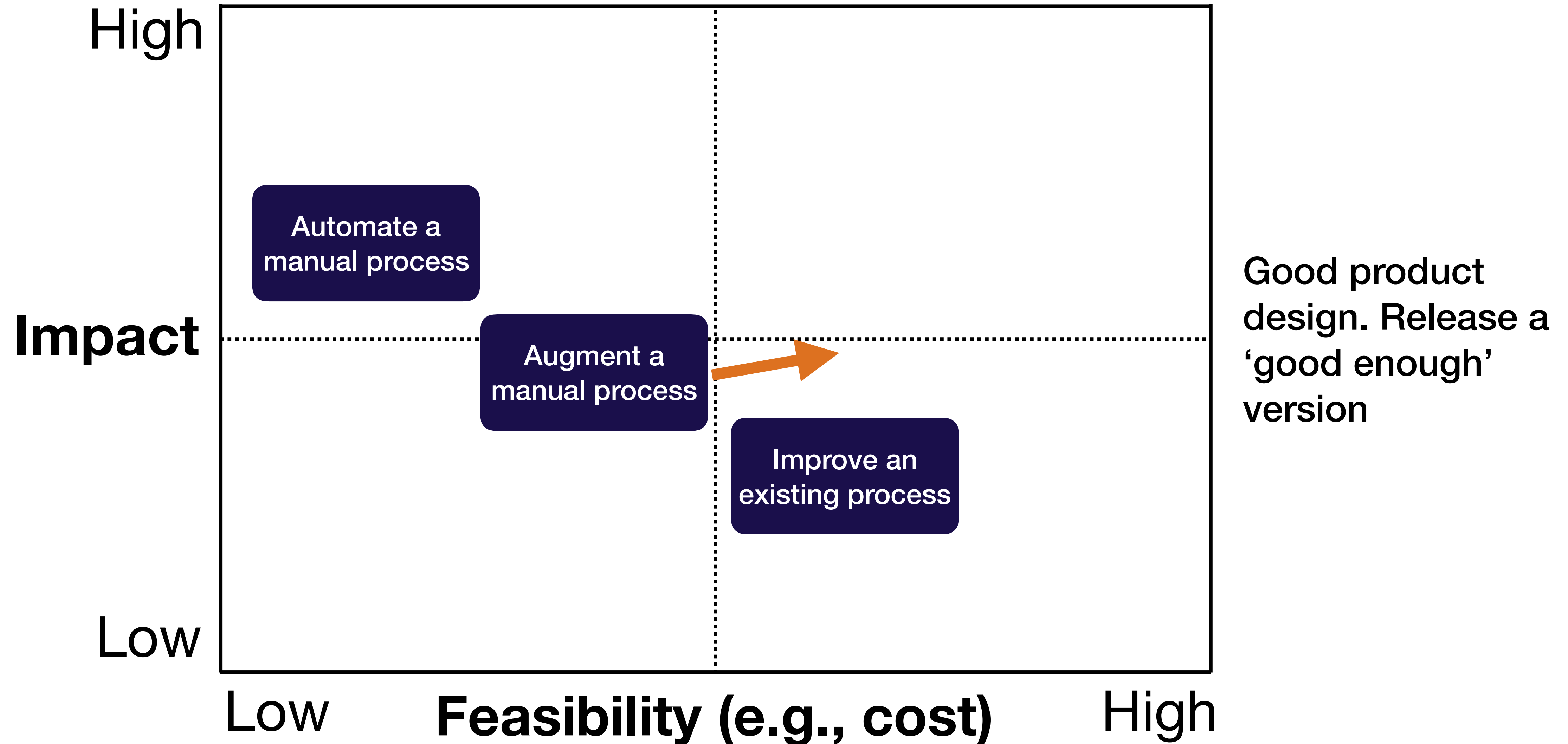
Machine learning project archetypes



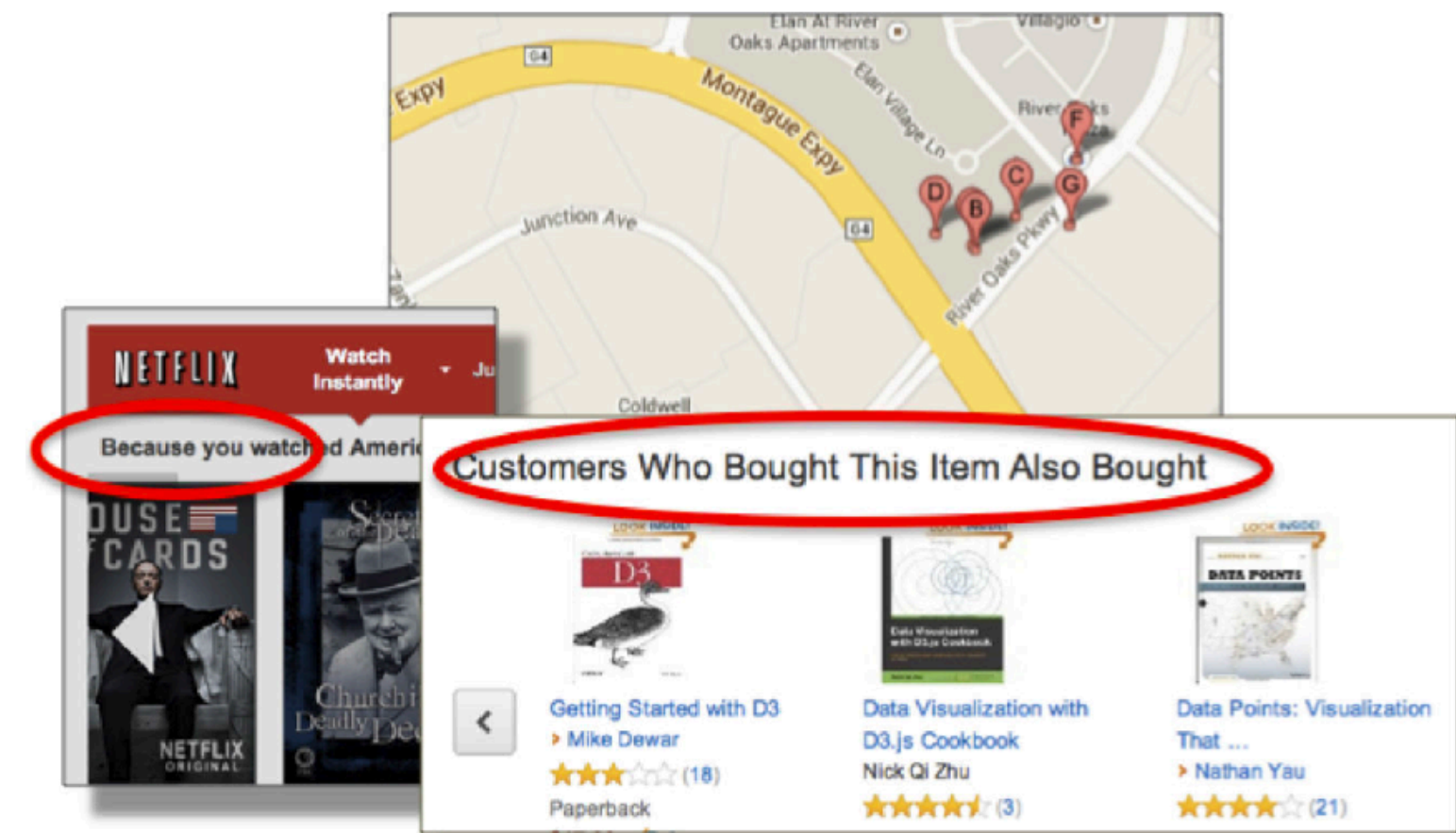
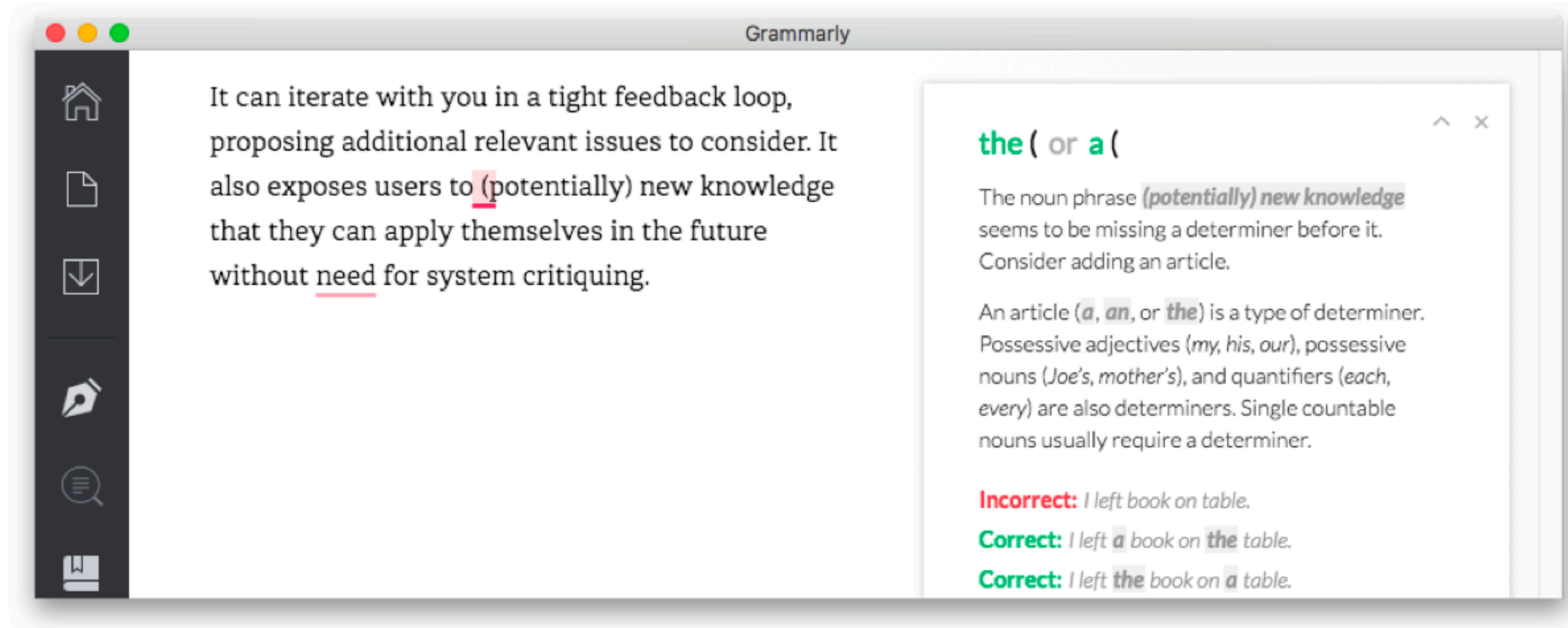
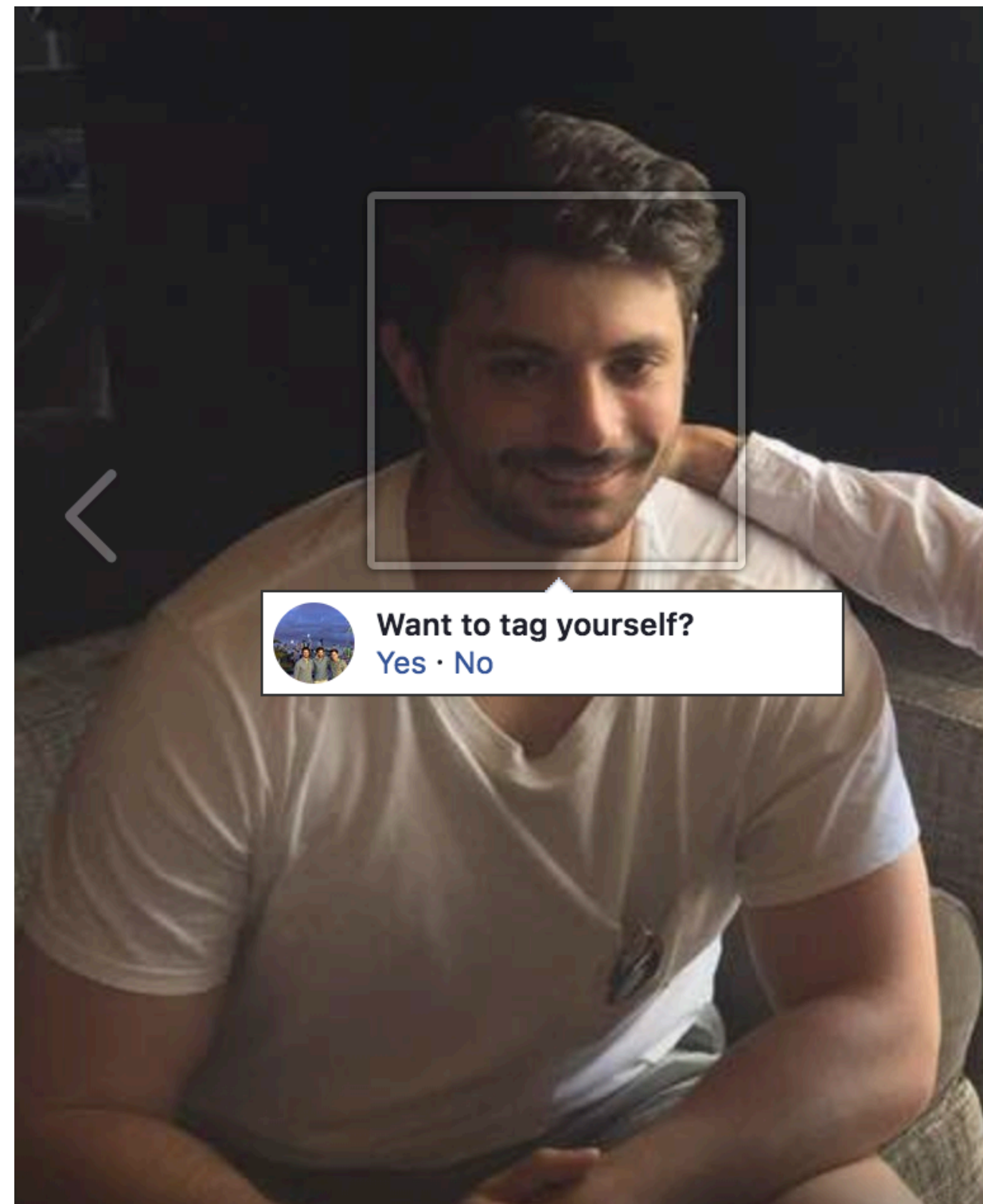
Machine learning project archetypes



Machine learning project archetypes

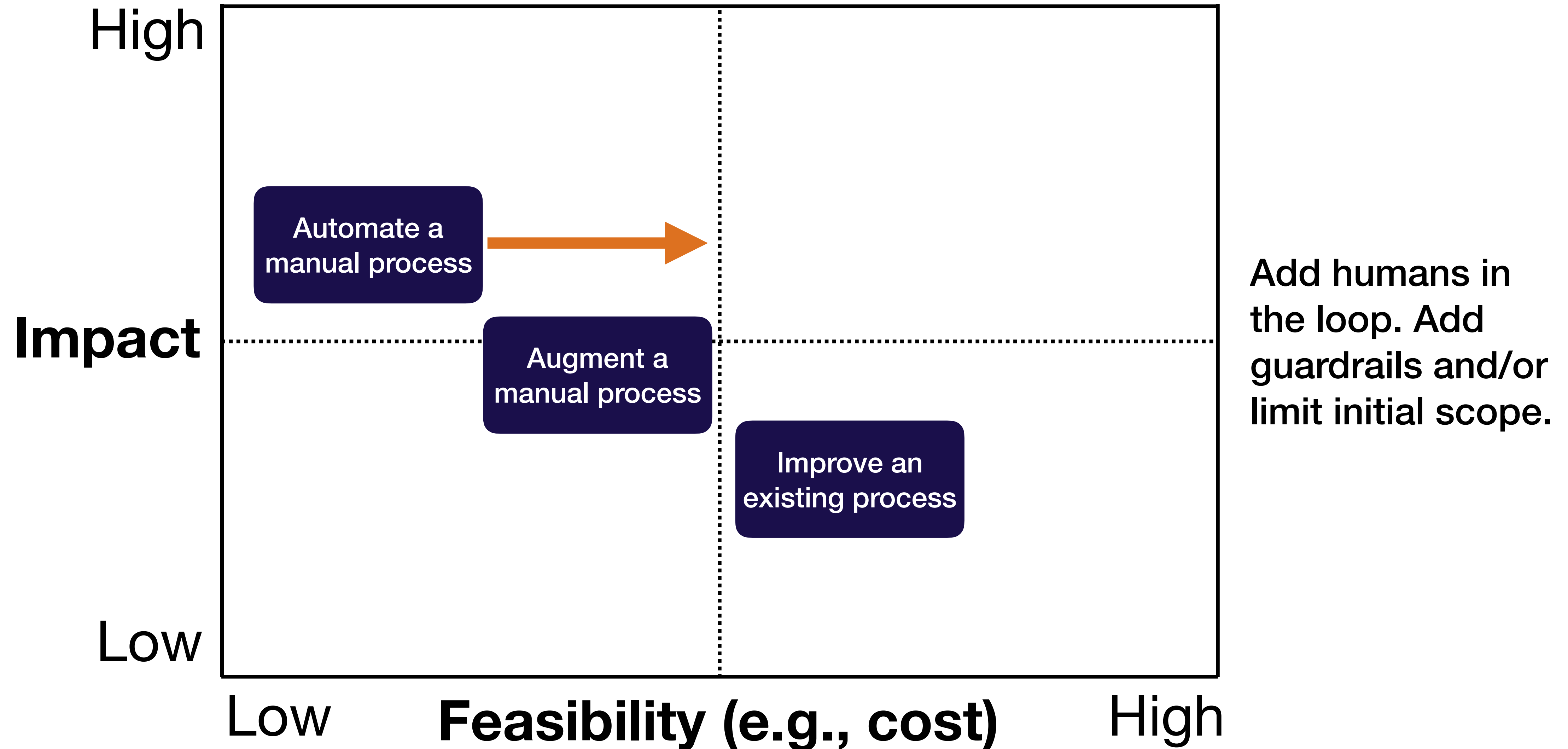


Product design can reduce need for accuracy



See "Designing Collaborative AI" (Ben Reinhardt and Belmer Negrillo):
https://medium.com/@Ben_Reinhardt/designing-collaborative-ai-5c1e8dbc8810






Machine learning project archetypes



Questions?



Module overview

-  Lifecycle
 - How to think about all of the activities in an ML project
-  Prioritizing projects
 - Assessing the feasibility and impact of your projects
-  Archetypes
 - The main categories of ML projects, and the implications for project management
-  **Metrics**
 - **How to pick a single number to optimize**
-  Baselines
 - How to know if your model is performing well

Key points for choosing a metric

- A. The real world is messy; you usually care about lots of metrics
- B. However, ML systems work best when optimizing a single number
- C. As a result, you need to pick a formula for combining metrics
- D. This formula can and will change!

Review of accuracy, precision, and recall

Confusion matrix

n=100	Predicted: NO	Predicted: YES	
Actual: NO	5	5	10
Actual: YES	45	45	90
	50	50	

Review of accuracy, precision, and recall

Confusion matrix

n=100	Predicted: NO	Predicted: YES	
Actual: NO	5	5	10
Actual: YES	45	45	90
	50	50	

$$\text{Accuracy} = \frac{\text{Correct}}{\text{Total}}$$

50%

Review of accuracy, precision, and recall

Confusion matrix

n=100	Predicted: NO	Predicted: YES	
Actual: NO	5	5	10
Actual: YES	45	45	90
	50	50	

Precision $\frac{\text{true positives}}{\text{true positives} + \text{false positives}}$
90%

Review of accuracy, precision, and recall

Confusion matrix

n=100	Predicted: NO	Predicted: YES	
Actual: NO	5	5	10
Actual: YES	45	45	90
	50	50	

Recall

true positives

Actual YES

50%

Why choose a single metric?

	Precision	Recall
Model 1	0.9	0.5
Model 2	0.8	0.7
Model 3	0.7	0.9

Which is best?

How to combine metrics

- Simple average / weighted average

Combining precision and recall

	Precision	Recall
Model 1	0.9	0.5
Model 2	0.8	0.7
Model 3	0.7	0.9

Combining precision and recall

	Precision	Recall	$(p + r) / 2$
Model 1	0.9	0.5	0.7
Model 2	0.8	0.7	0.75
Model 3	0.7	0.9	0.8

Combining precision and recall

	Precision	Recall	$(p + r) / 2$
Model 1	0.9	0.5	0.7
Model 2	0.8	0.7	0.75
Model 3	0.7	0.9	0.8

How to combine metrics

- Simple average / weighted average

How to combine metrics

- Simple average / weighted average
- Threshold $n-1$ metrics, evaluate the n th

Thresholding metrics

Choosing which metrics to threshold

- Domain judgment (e.g., which metrics can you engineer around?)
- Which metrics are least sensitive to model choice?
- Which metrics are closest to desirable values?

Choosing threshold values

- Domain judgment (e.g., what is an acceptable tolerance downstream? What performance is achievable?)
- How well does the baseline model do?
- How important is this metric right now?

Combining precision and recall

	Precision	Recall	$(p + r) / 2$
Model 1	0.9	0.5	0.7
Model 2	0.8	0.7	0.75
Model 3	0.7	0.9	0.8

Combining precision and recall

	Precision	Recall	$(p + r) / 2$	$p @ (r > 0.6)$
Model 1	0.9	0.5	0.7	0.0
Model 2	0.8	0.7	0.75	0.8
Model 3	0.7	0.9	0.8	0.7

Combining precision and recall

	Precision	Recall	$(p + r) / 2$	$p @ (r > 0.6)$
Model 1	0.9	0.5	0.7	0.0
Model 2	0.8	0.7	0.75	0.8
Model 3	0.7	0.9	0.8	0.7

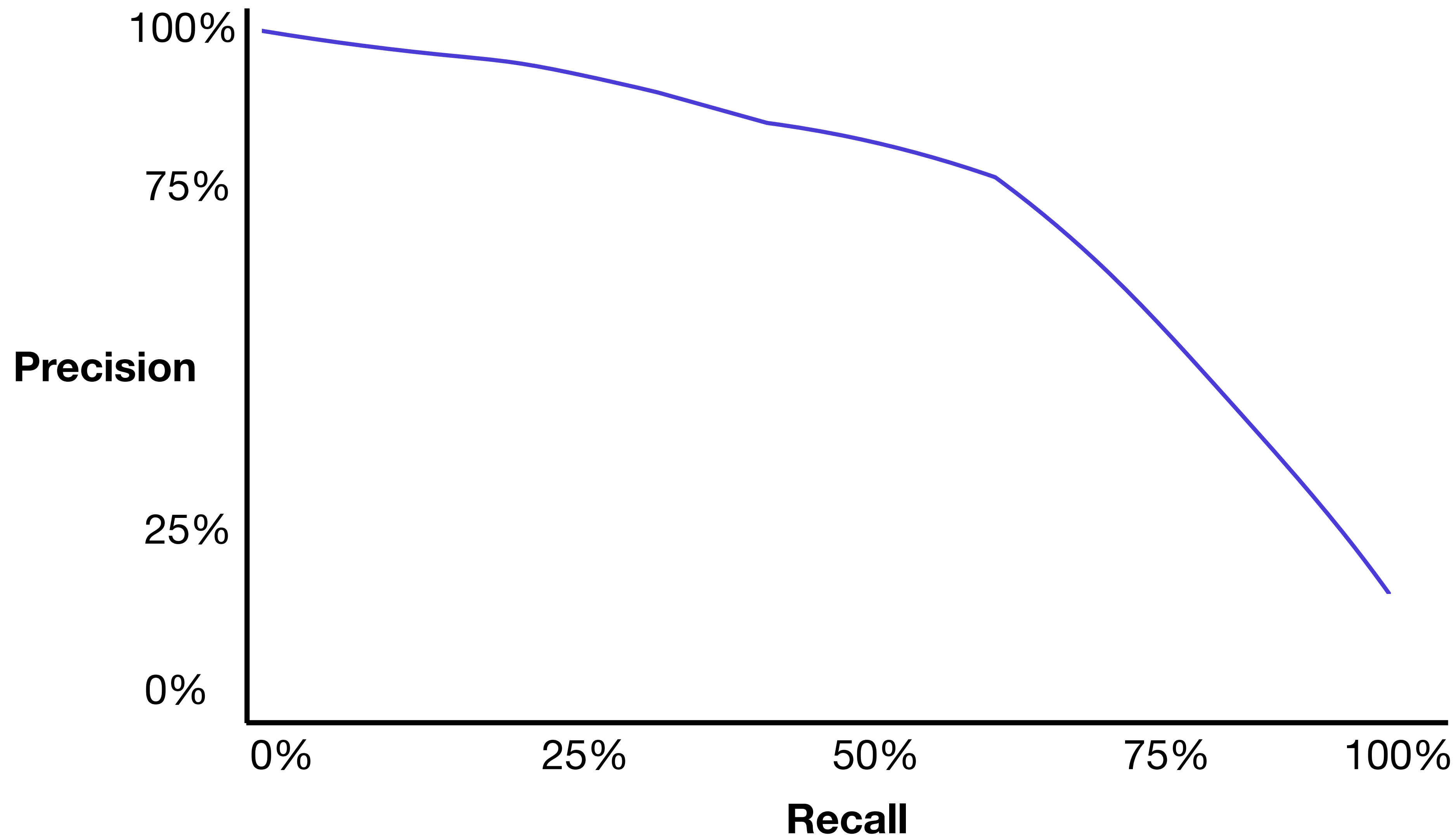
How to combine metrics

- Simple average / weighted average
- Threshold $n-1$ metrics, evaluate the n th

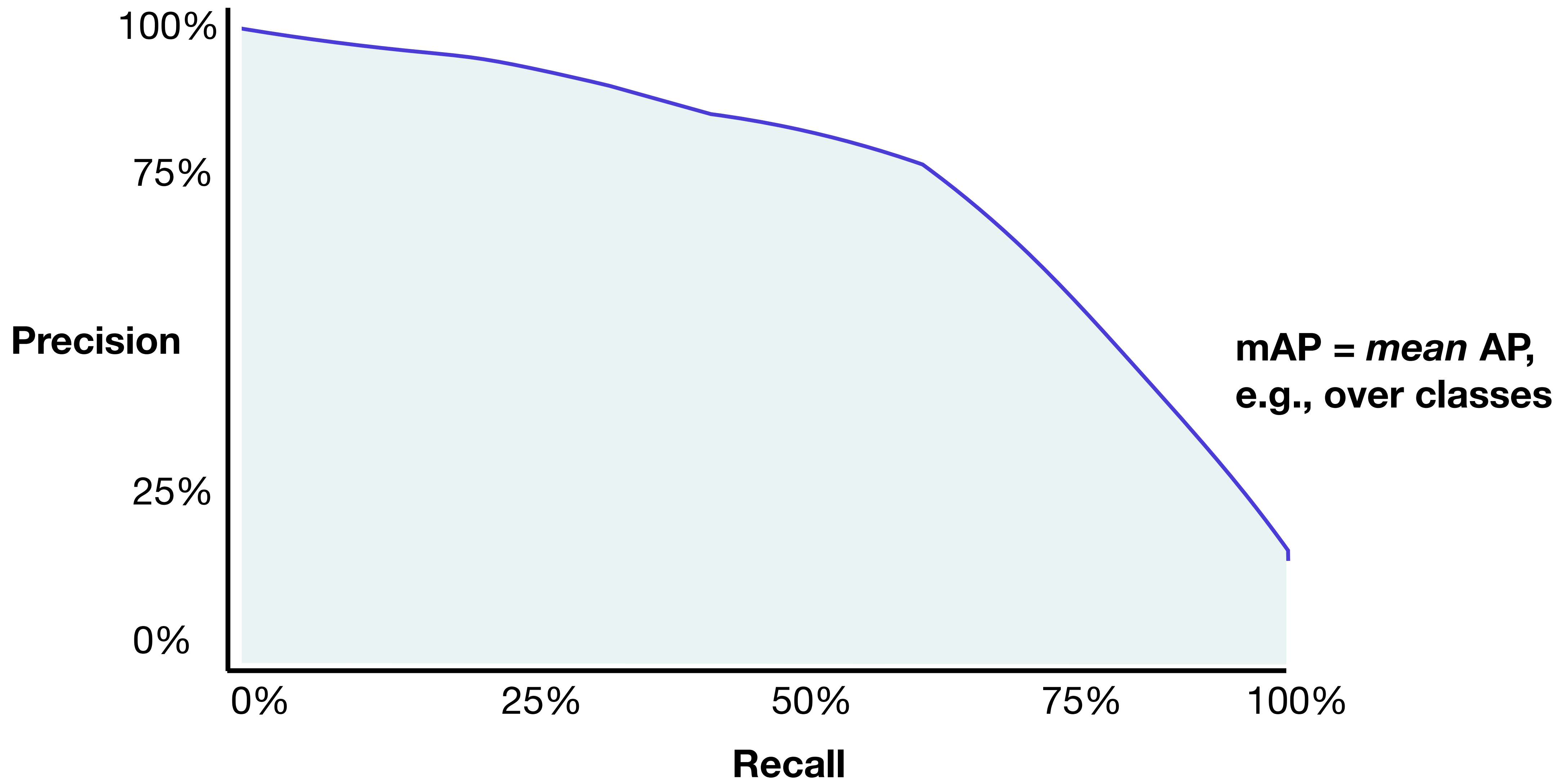
How to combine metrics

- Simple average / weighted average
- Threshold $n-1$ metrics, evaluate the n th
- More complex / domain-specific formula

Domain-specific metrics: mAP



Domain-specific metrics: mAP



Combining precision and recall

	Precision	Recall	$(p + r) / 2$	$p @ (r > 0.6)$
Model 1	0.9	0.5	0.7	0.0
Model 2	0.8	0.7	0.75	0.8
Model 3	0.7	0.9	0.8	0.7

Combining precision and recall

	Precision	Recall	$(p + r) / 2$	$p @ (r > 0.6)$	mAP
Model 1	0.9	0.5	0.7	0.0	0.7
Model 2	0.8	0.7	0.75	0.8	0.6
Model 3	0.7	0.9	0.8	0.7	0.6

Combining precision and recall

	Precision	Recall	$(p + r) / 2$	$p @ (r > 0.6)$	mAP
Model 1	0.9	0.5	0.7	0.0	0.7
Model 2	0.8	0.7	0.75	0.8	0.6
Model 3	0.7	0.9	0.8	0.7	0.6

Example: choosing a metric for pose estimation



(x, y, z) **Position (L2 loss)**

(ϕ, θ, ψ) **Orientation (L2 loss)**

t **Prediction time**

Xiang, Yu, et al. "PoseCNN: A Convolutional Neural Network for 6D Object Pose Estimation in Cluttered Scenes." *arXiv preprint arXiv:1711.00199* (2017).

Example: choosing a metric for pose estimation

- **Enumerate requirements**
 - Downstream goal is real-time robotic grasping
 - Position error must be $<1\text{cm}$, not sure exactly how precise is needed
 - Angular error <5 degrees
 - Must run in 100ms to work in real-time

Example: choosing a metric for pose estimation

- Enumerate requirements
- **Evaluate current performance**
- Train a few models

Example: choosing a metric for pose estimation

- Enumerate requirements
- Evaluate current performance
- **Compare current performance to requirements**
 - Position error between 0.75 and 1.25cm (depending on hyperparameters)
 - All angular errors around 60 degrees
 - Inference time ~300ms

Example: choosing a metric for pose estimation

- Enumerate requirements
- Evaluate current performance
- **Compare current performance to requirements**
 - Prioritize angular error
 - Threshold position error at 1cm
 - Ignore run time for now

Example: choosing a metric for pose estimation

- Enumerate requirements
- Evaluate current performance
- Compare current performance to requirements
- **Revisit metric as your numbers improve**




Key points for choosing a metric

- A. The real world is messy; you usually care about lots of metrics
- B. However, ML systems work best when optimizing a single number
- C. As a result, you need to pick a formula for combining metrics
- D. This formula can and will change!

Questions?



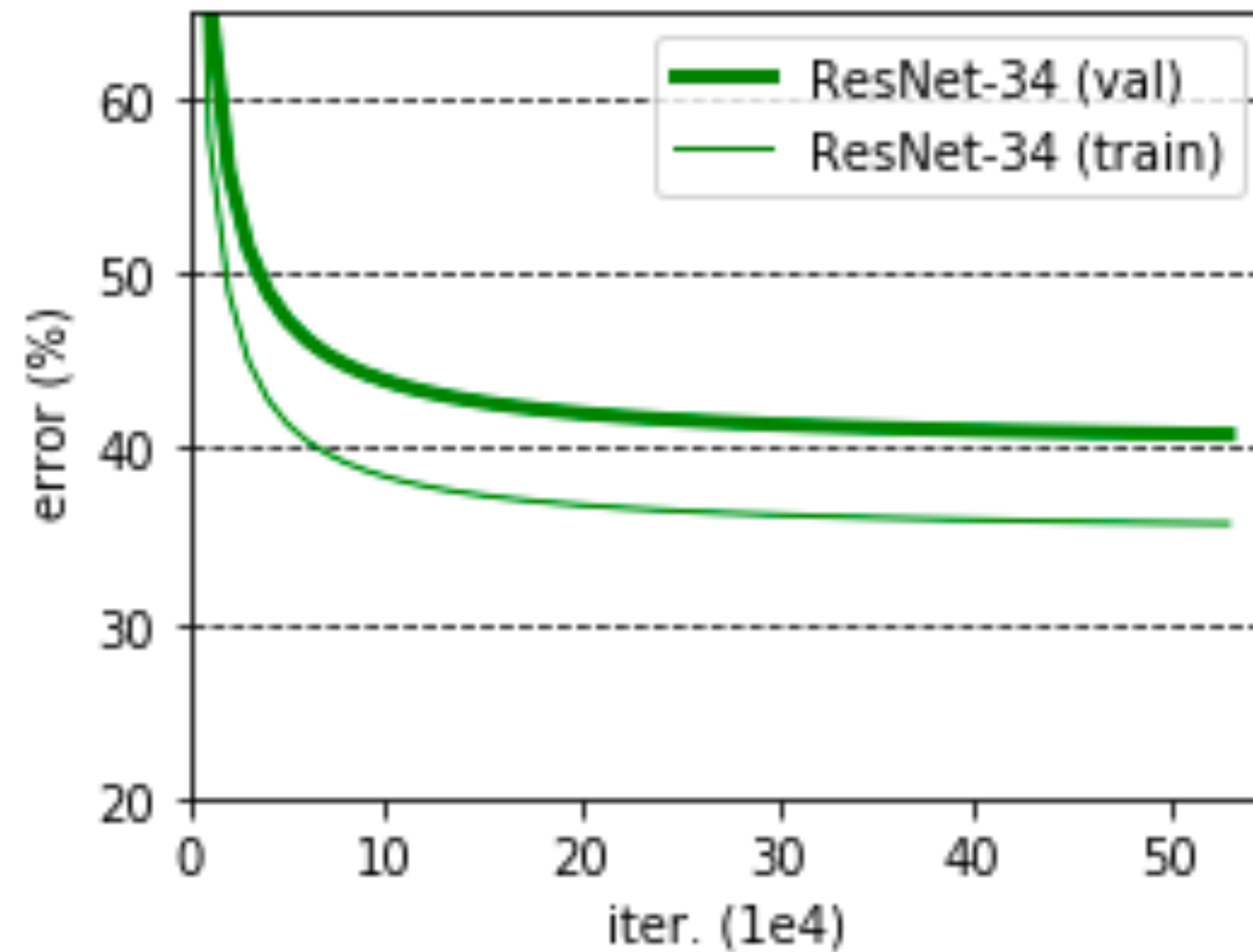
Module overview

-  Lifecycle
 - How to think about all of the activities in an ML project
-  Prioritizing projects
 - Assessing the feasibility and impact of your projects
-  Archetypes
 - The main categories of ML projects, and the implications for project management
-  Metrics
 - How to pick a single number to optimize
-  **Baselines**
 - **How to know if your model is performing well**

Key points for choosing baselines

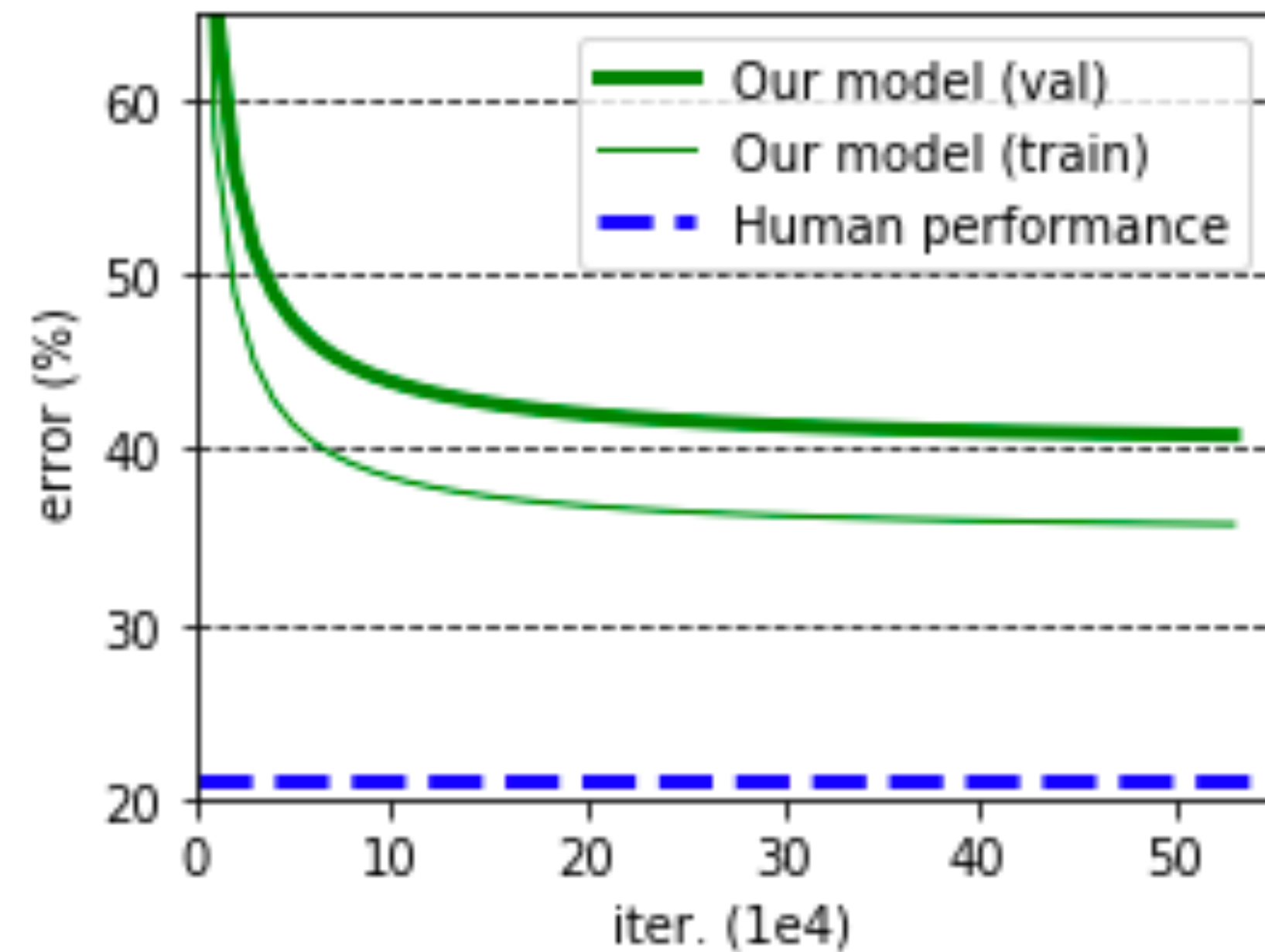
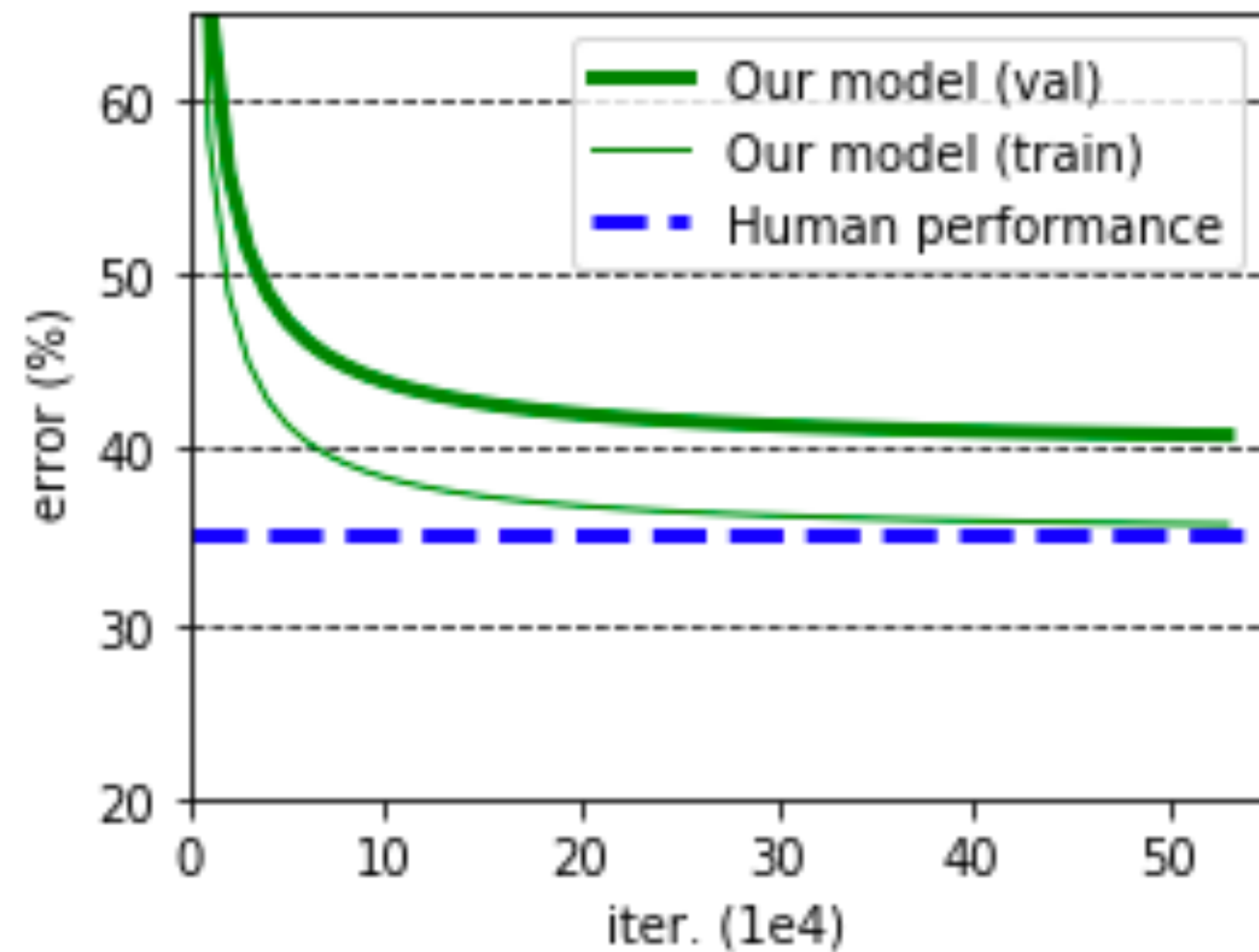
- A. Baselines give you a lower bound on expected model performance
- B. The tighter the lower bound, the more useful the baseline (e.g., published results, carefully tuned pipelines, & human baselines are better)

Why are baselines important?



Why are baselines important?

Same model, different baseline → different next steps



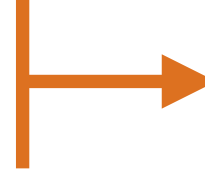
Where to look for baselines

External
baselines

- Business / engineering requirements

Where to look for baselines

External
baselines

- Business / engineering requirements
- Published results  **Make sure comparison is fair!**

Where to look for baselines

External baselines

- Business / engineering requirements
- Published results

Internal baselines

- Scripted baselines (e.g., OpenCV, rules-based)

Where to look for baselines

External baselines

- Business / engineering requirements
- Published results

Internal baselines

- Scripted baselines (e.g., OpenCV, rules-based)
- Simple ML baselines (e.g., bag of words, linear regression)

Where to look for baselines

External baselines

- Business / engineering requirements
 - Published results
-

Internal baselines

- Scripted baselines (e.g., OpenCV, rules-based)
- Simple ML baselines (e.g., bag of words, linear regression)
- Human performance

How to create good human baselines

Quality of baseline

Low

Random people (e.g., Amazon Turk)

Ensemble of random people

Domain experts (e.g., doctors)

Deep domain experts (e.g., specialists)

Mixture of experts

Ease of data collection

High

Low

How to create good human baselines

- Highest quality that allows more data to be labeled easily
- More specialized domains need more skilled labelers
- Find cases where model performs worse and concentrate data collection there

More on labeling in data lecture!

Key points for choosing baselines

- A. Baselines give you a lower bound on expected model performance
- B. The tighter the lower bound, the more useful the baseline (e.g., published results, carefully tuned pipelines, human baselines are better)

Questions?

Conclusion



- **ML projects are iterative. Deploy something fast to begin the cycle.**



- **Choose projects that that are high impact with low cost of wrong predictions**



- **The secret sauce to making projects work well is to build automated data flywheels**



- **In the real world you care about many things, but you should always have just one you're working on**



- **Good baselines help you invest your effort the right way**

Where to go to learn more

- Andrew Ng's "Machine Learning Yearning"
- Andrej Karpathy's "Software 2.0"
- Agrawal's "The Economics of AI"

Thank you!