

# An introduction to contextual bandits

Wilson Pok

# Goal and rules of the game

**Sequential decision-making  
under uncertainty**

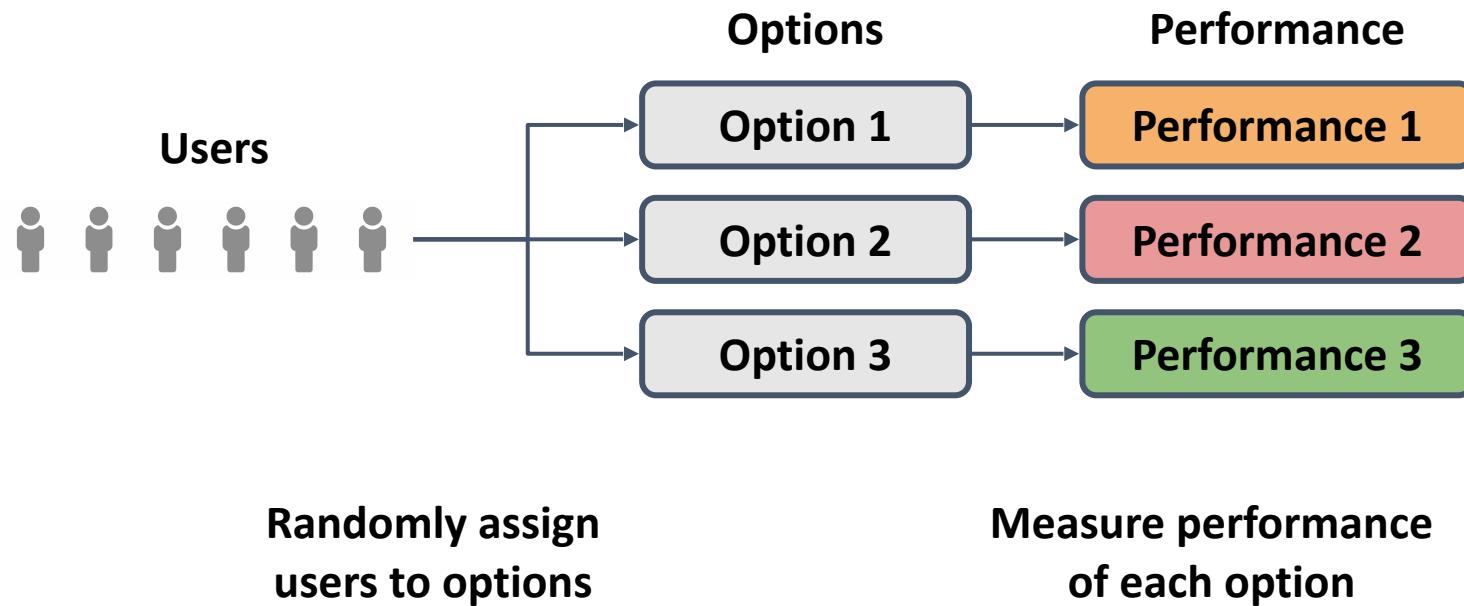
Each round is independent

Learn by trying - but it costs to try

Only see results of decisions made

*How do you quickly and confidently learn the best decision each round?*

# We know how experimentation works

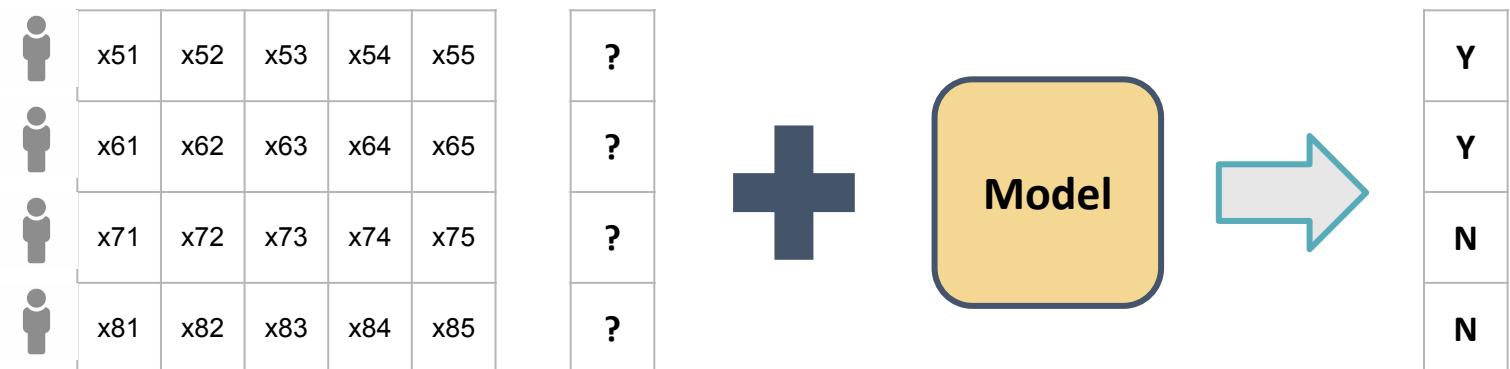


# We know how machine learning works

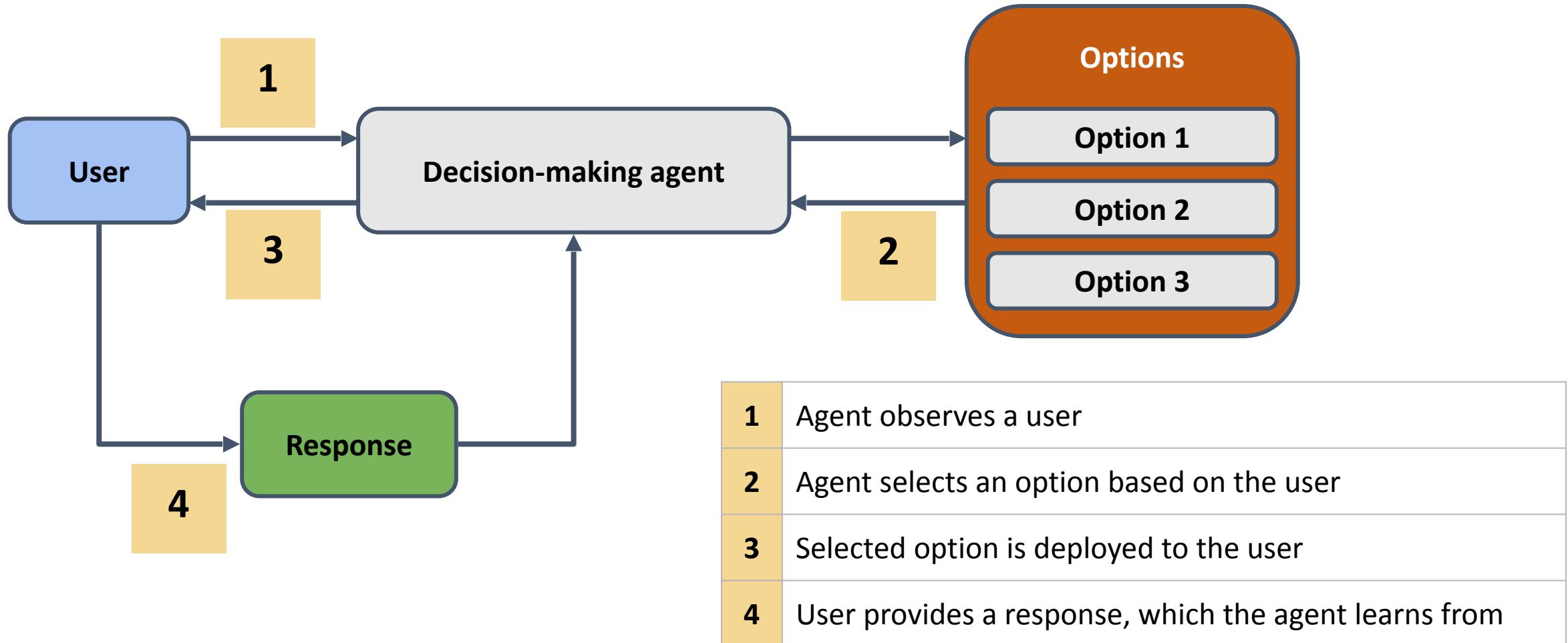
Create model from  
training data



Use model to  
predict new data

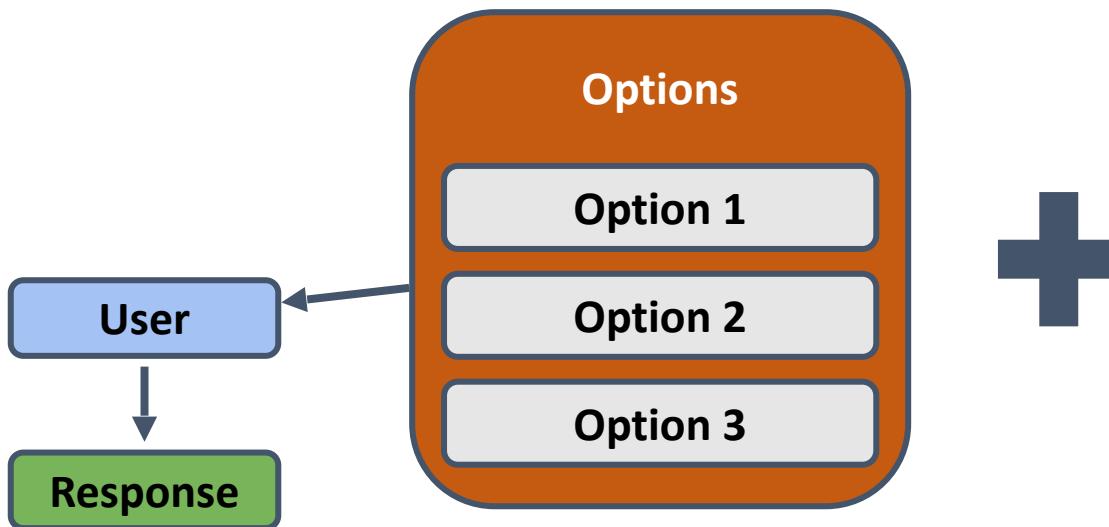


# How real-time decision-making works

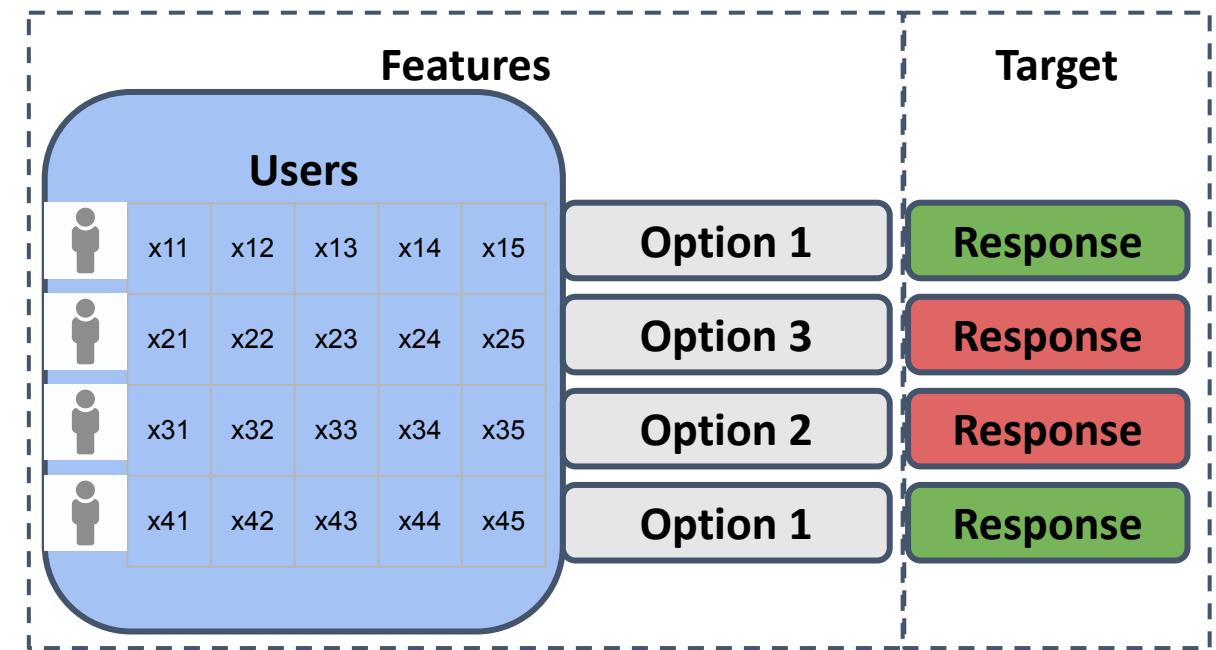


# This combines experimentation with ML

## Experimentation



## Machine Learning



# Explore/exploit tradeoff

## Experimentation

On one hand, we need experimentation to *explore* our options.

By trying all options on different types of users, we can generate enough training data confidently predict their responses.



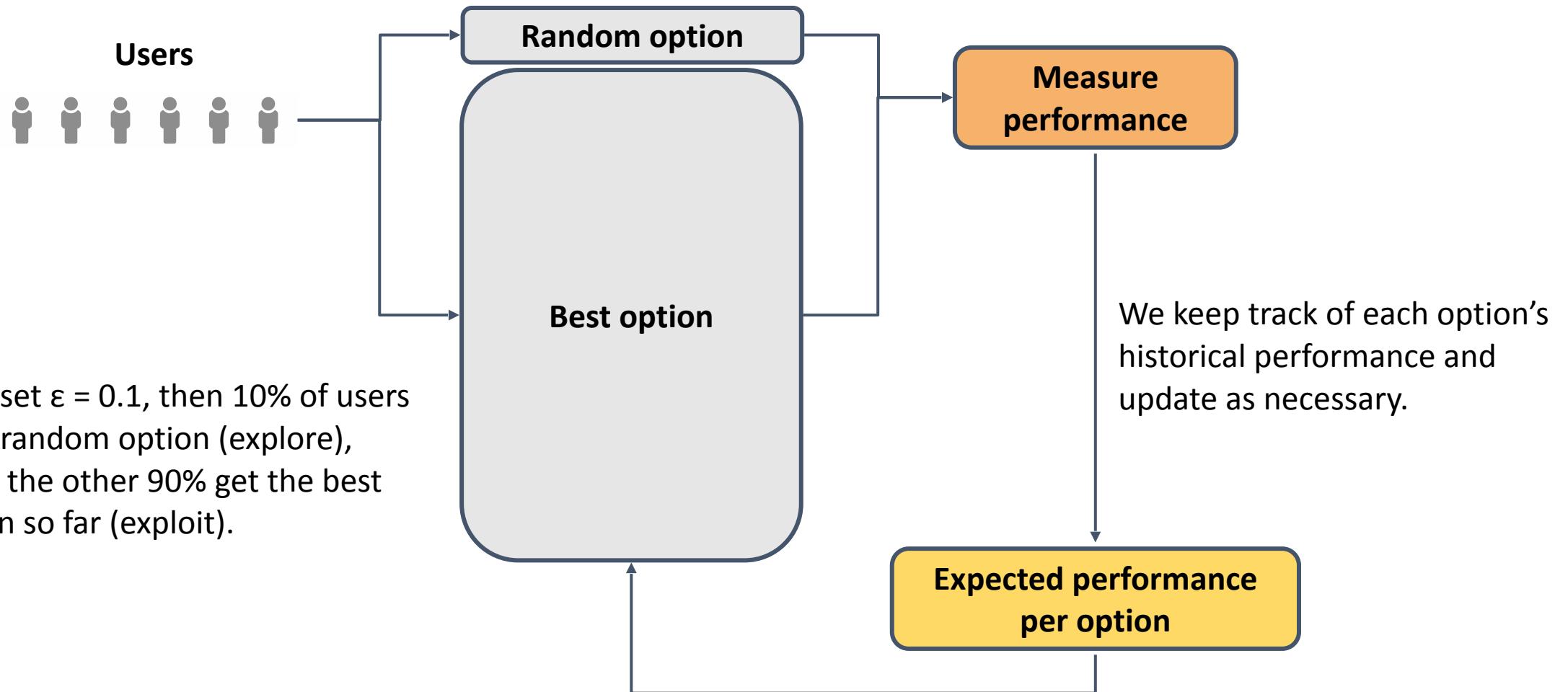
## Machine Learning

On the other hand, we need to use the machine learning model to *exploit* what we know.

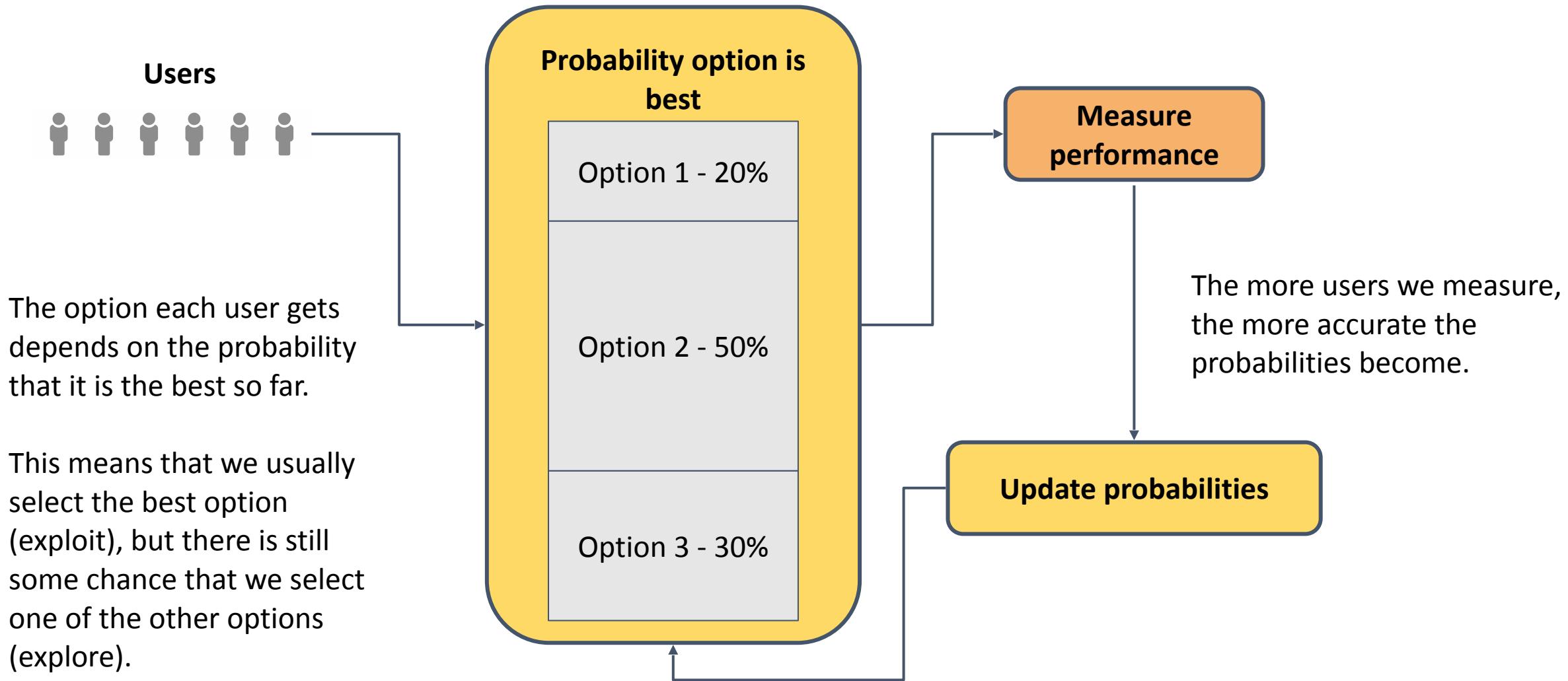
If we are confident that a particular option is best for a certain type of user, then why would we give that user anything else?

If only there were some way to tune the amount the decision-making agent explores/exploits...

# A simple approach is $\epsilon$ -greedy



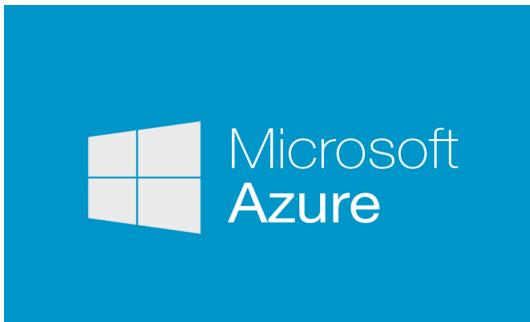
# Another approach is Thompson sampling



# Practical applications

<b>Recommender systems</b>	Recommending news articles to readers on a news website, personalising marketing content
<b>Healthcare</b>	<ul style="list-style-type: none"><li>- Adaptive allocation in clinical trials - give more patients treatments that look more promising</li><li>- Dosing strategies - finding the right dose for drugs that are highly variable depending on the patient</li></ul>
<b>Finance</b>	Automated strategies for constructing portfolios that balance risk and reward
<b>Dynamic pricing</b>	Automated strategies for maximising revenue based on observed demand
<b>Influence maximisation</b>	Maximizing the number of users that become aware of a product by selecting a set of seed users to expose the product to. The agent finds the best influencers in a social network by repeatedly interacting with it.
<b>Dialogue systems</b>	Creating conversational agents that can learn what the best response is for a given input

# Current implementations



[Personalizer](#)



Adobe Target  
[Automated Personalization](#)



[Personalization](#)

# Etymology

# Multi-armed bandits

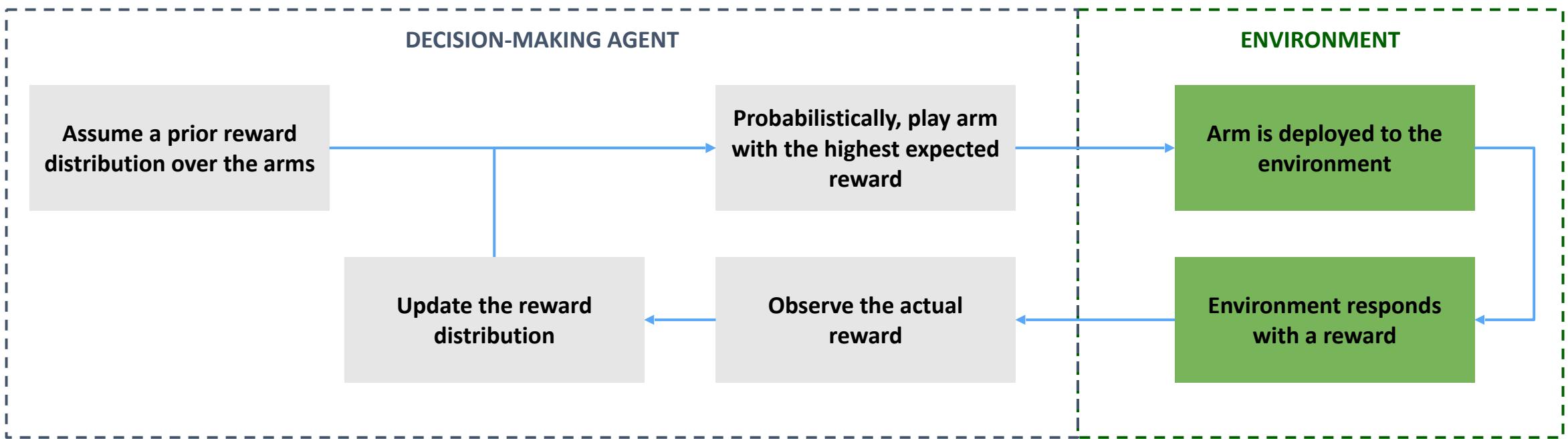
You are given a row of slot machines, each with different payout rates. The only way to estimate the payout rate of a slot machine is to play it.

*How do you maximise your winnings?*

Specifically, how do you identify the best arm to play **quickly and confidently**?

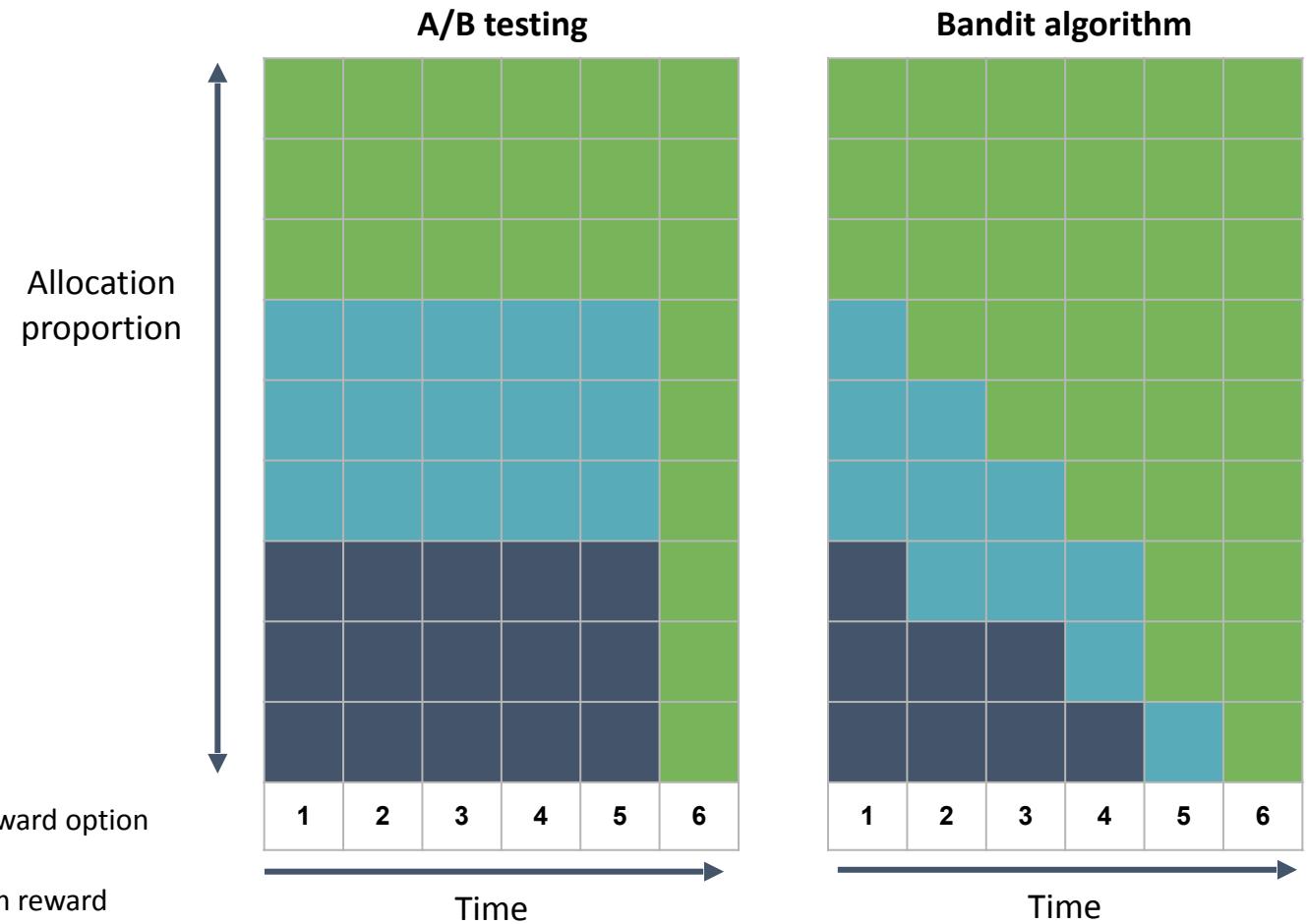


# Learn by interacting with environment

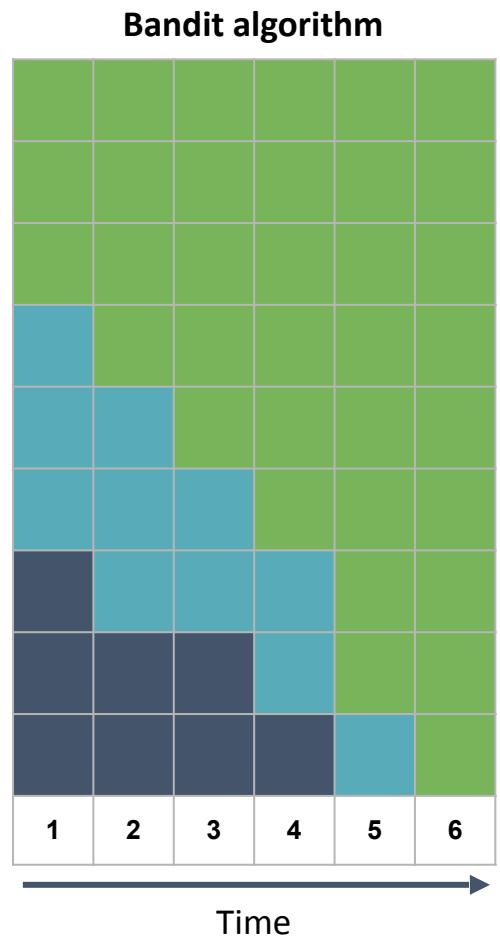


# Bandits automate A/B testing

A/B testing	Bandit algorithms
Need to wait for a final answer (losing reward until then)	Converges to best option throughout run (gaining reward continuously)
Needs to be rerun if environment changes	Can be left running, will adapt to environment automatically
Can only find the best overall option	Can do personalisation ( <b>contextual bandits</b> )



With A/B testing, the test must be run completely to find the option with the highest reward.



In contrast, bandit algorithms can continuously shift the allocation toward high performing options faster.

# Contextual bandits

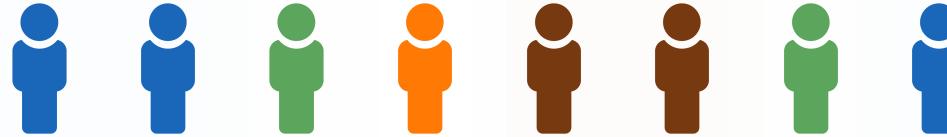
Now suppose that payout rates change depending on the attributes of the player (i.e. rewards depend on *context*). We have a series of players (with different attributes) who will play one arm each, sequentially.

***How do we maximise the overall winnings?***

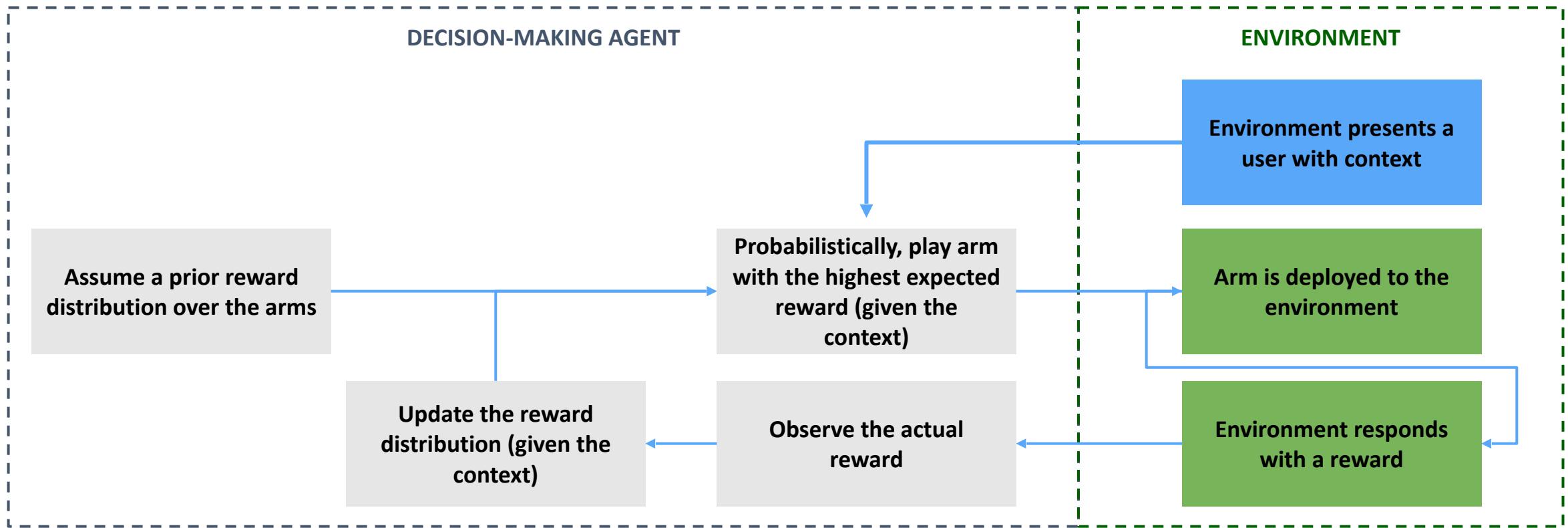
## ***Why do we care about this?***

Solving this problem allows us to tailor the arm to fit each particular context.

Specifically, contextual bandits enable personalised actions to be applied to individual users.



# Learn by interacting with environment (CB)



# Multi-armed vs contextual bandits

Multi-armed bandits	Contextual bandits
Automated A/B testing	Automated personalisation for each user
Decisions are made based on past performance of each option	Decisions are made based on past performance of each option on each user type
Aims to give most users the best overall option	Aims to give most users the best option for them