

Reinforcement Learning for the real world

Harshad Khadilkar
Tata Consultancy Services Ltd.



Motivation: Why RL?

This is about letting an ecosystem of machines teach itself superhuman capabilities

Why?

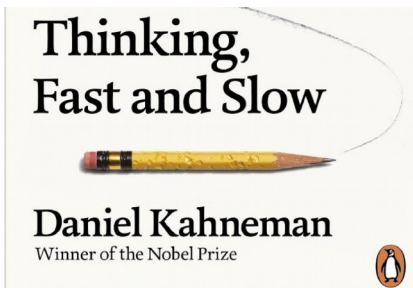
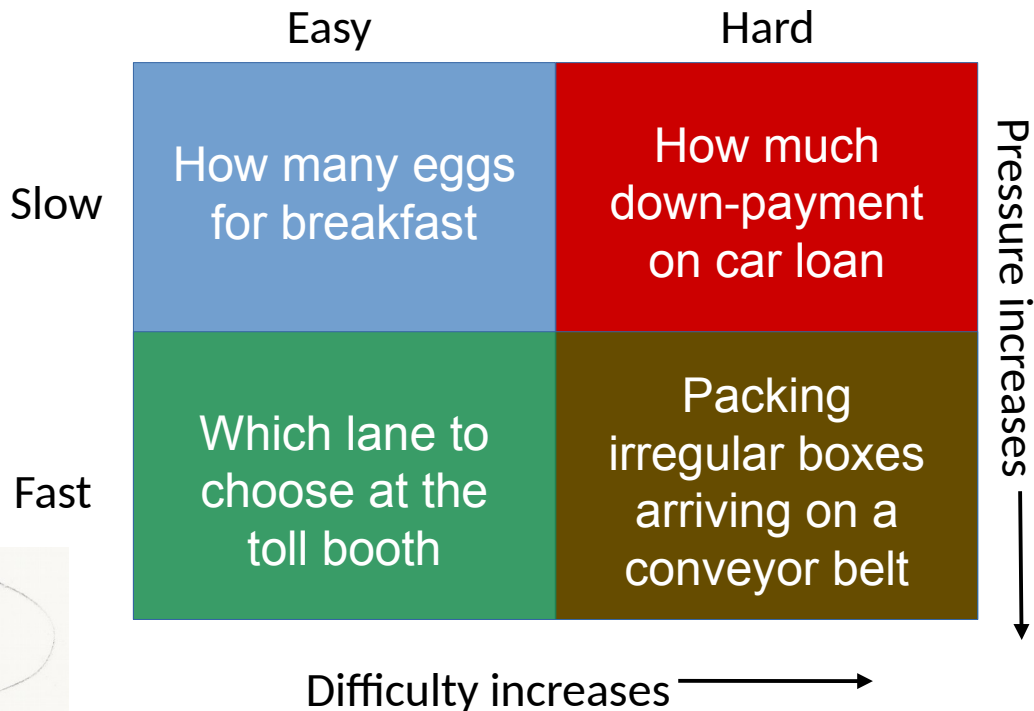
“Because it’s there”

- George Mallory (1923), when asked why he wanted to climb Mt. Everest



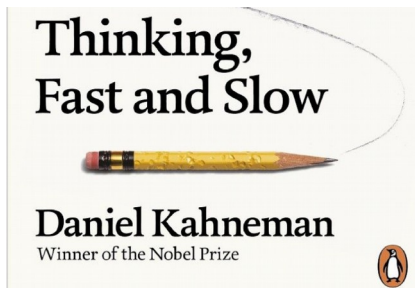
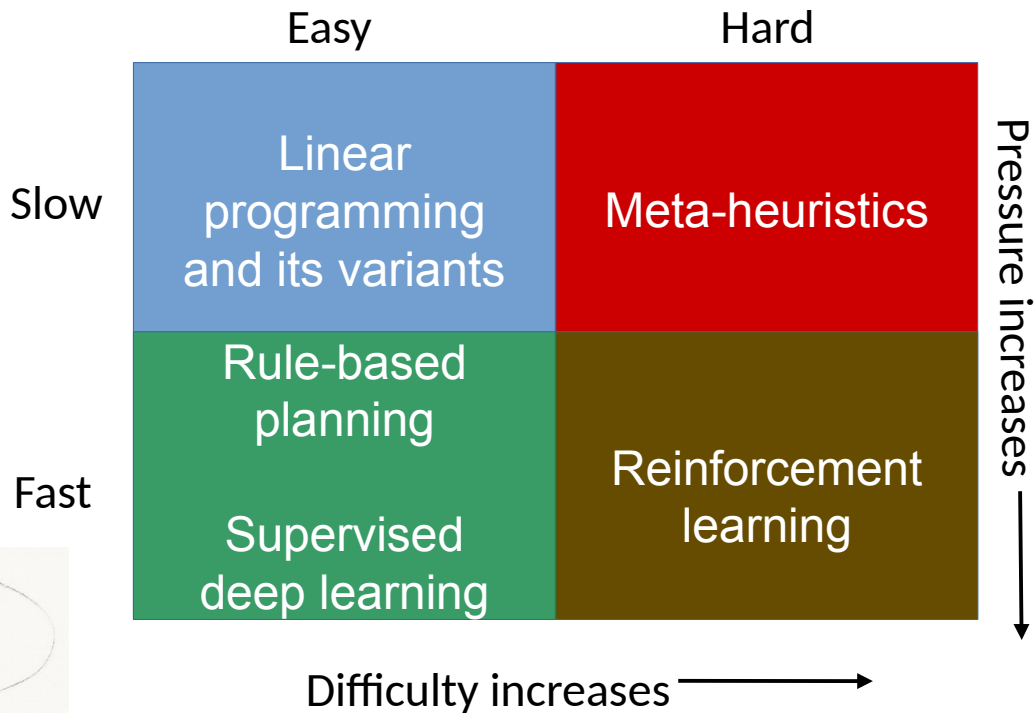
2 Motivation

RL in the optimization space



2 Motivation

RL in the optimization space



Necessary conditions: Answer YES to all of the following

Use for tasks that humans find hard to do (or to do well) → No ideal reference

When time is short → Can't search or solve in real-time

When the system is hard to define, or complex → No analytical relationships

“The most important training in Unseen University [for wizards] wasn't how to do magic,
but to know when not to use it” - Terry Pratchett

This is about letting an ecosystem of machines teach itself superhuman capabilities

Why?

“Because it’s there”

- George Mallory (1923), when asked why he wanted to climb Mt. Everest



How?

Let the algorithm explore the environment on its own, while **learning from experience**

Reinforcement learning



Briefly:
What is it?

Learning to maximise long-term reward through interaction with the environment



3 How RL works

Context: Existing work

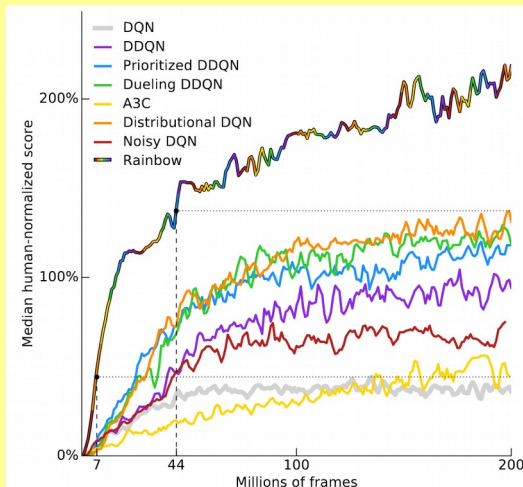
Is this a new idea?

Not at all.

Aerospace:
Adaptive control

Ops Research:
Dynamic programming

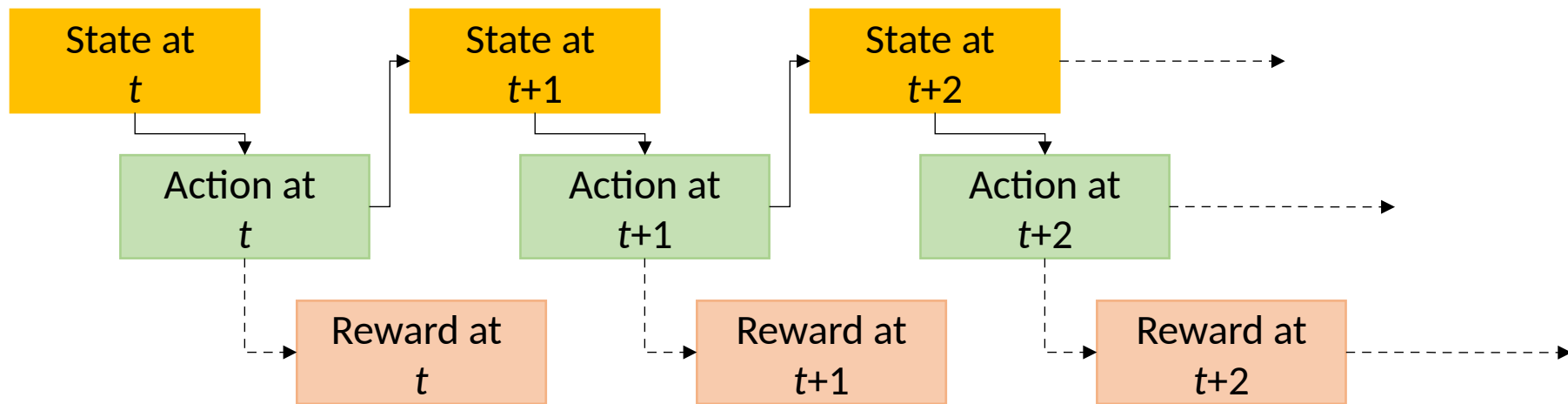
Computer science:
Reinforcement learning



3

How RL works

Anatomy of an RL problem

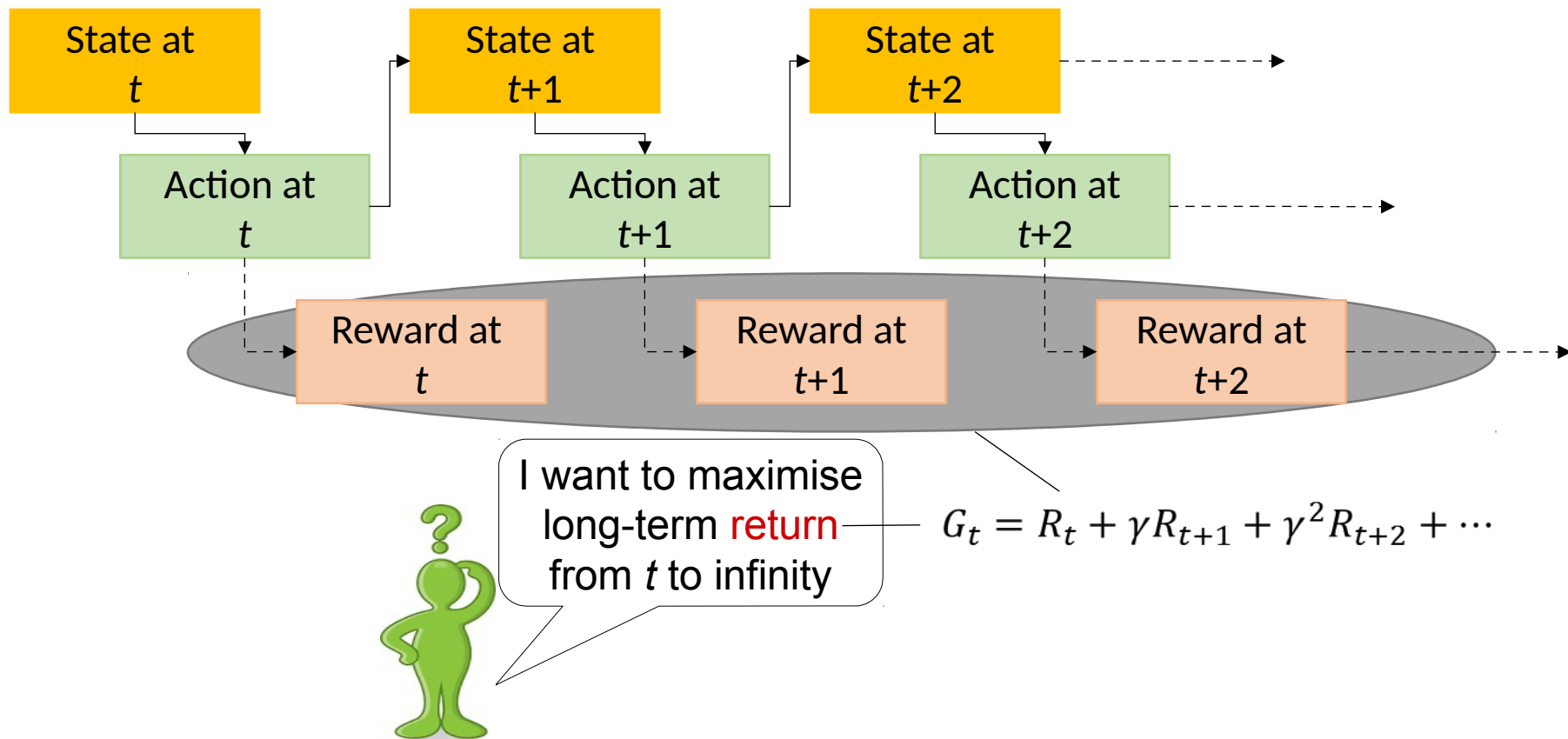


Strictly speaking, must be a **Markov Decision Process** defined by

(States, Actions, Rewards, Transitions, Discount factor)

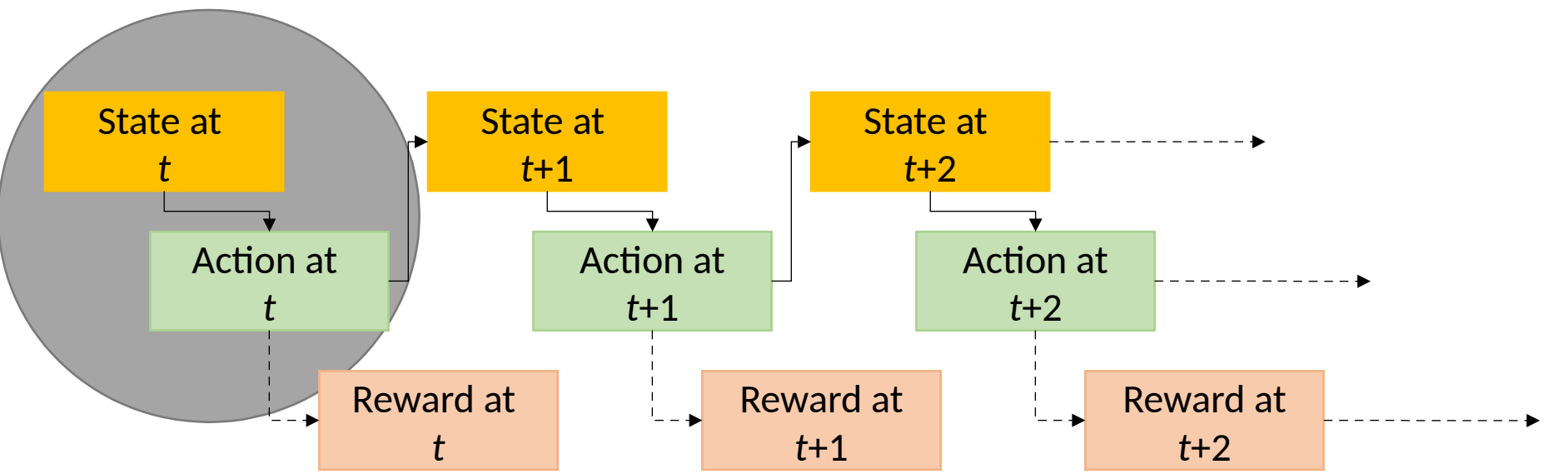
3 How RL works

Anatomy of an RL problem



3 How RL works

Anatomy of an RL problem



Use neural network, with RHS of equation providing labels

Value-based Deep RL

Can be any function from (state, action) \rightarrow scalar

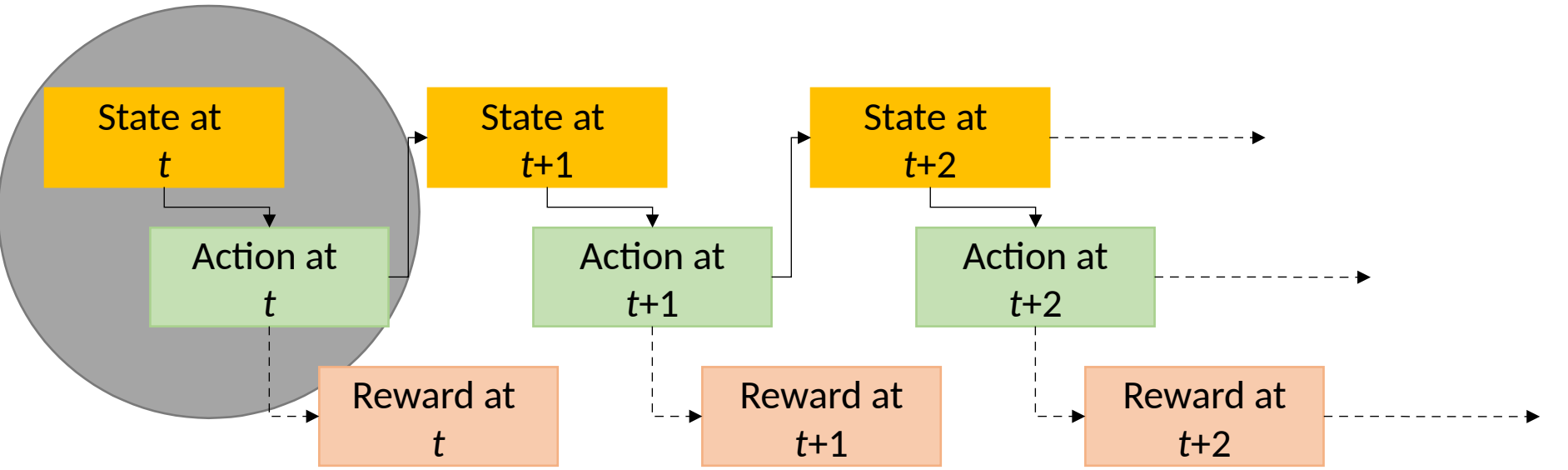
This is unknown at t , but is 'known' at $t+1$

All value based approaches

$$Q(s_t, a_t) = R_t + \gamma Q(s_{t+1}, a_{t+1})$$

3 How RL works

Anatomy of an RL problem



Use neural network,
with gradients driving
the training

Policy-based **Deep RL**

Gradient of expected
reward with respect to
each element of

Goal: Compute
parameters θ that
maximize reward

All **policy** based approaches

$\Pi_{\theta} : S_t \rightarrow a_t$

Alternative approach

The bad news
These ideas work brilliantly in games,
but not in real life

Why not?

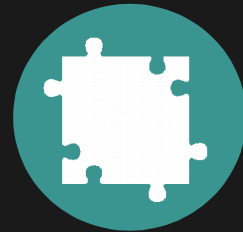


1. Large scale
2. Variable scale

3. Complexity

4. Limited compute

5. Explainability
requirement

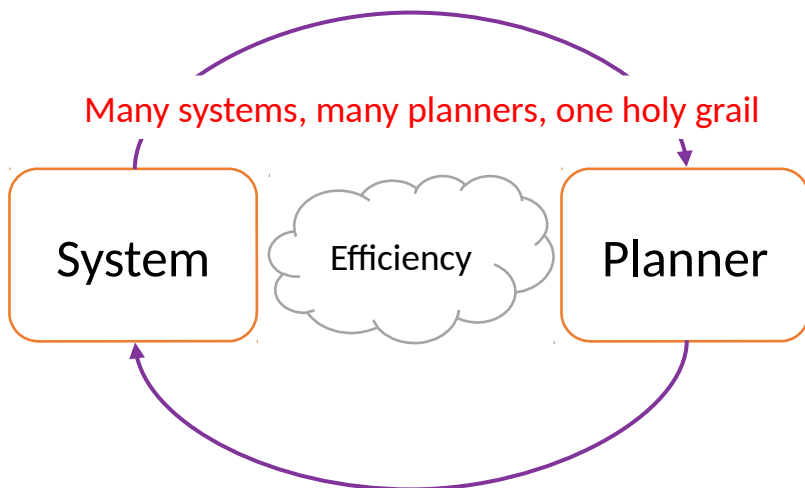


RL in the
real world

4 RL in the real world

One-slide summary of past work

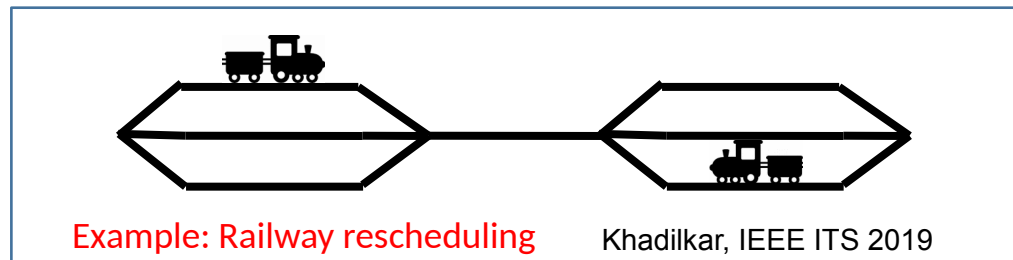
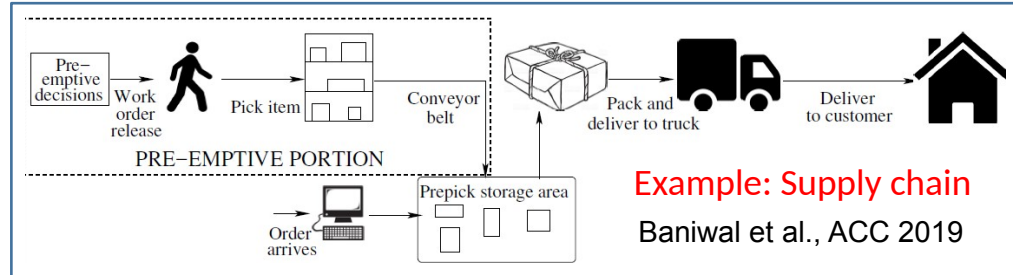
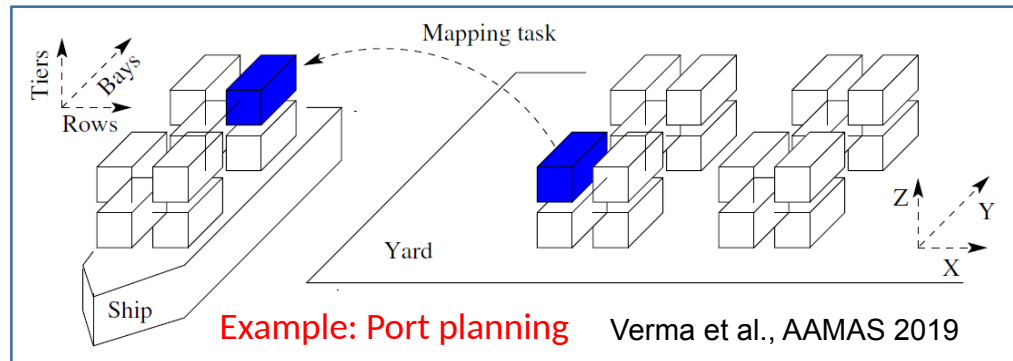
Many systems, many planners, one holy grail



Problem: Systems do not operate in silos ...
... but planners/controllers do

Goal: Build optimal planning & control algorithms that,

1. Operate in real-time (online)
2. Work without human-labelled historical data
3. Adapt automatically to changes in the environment



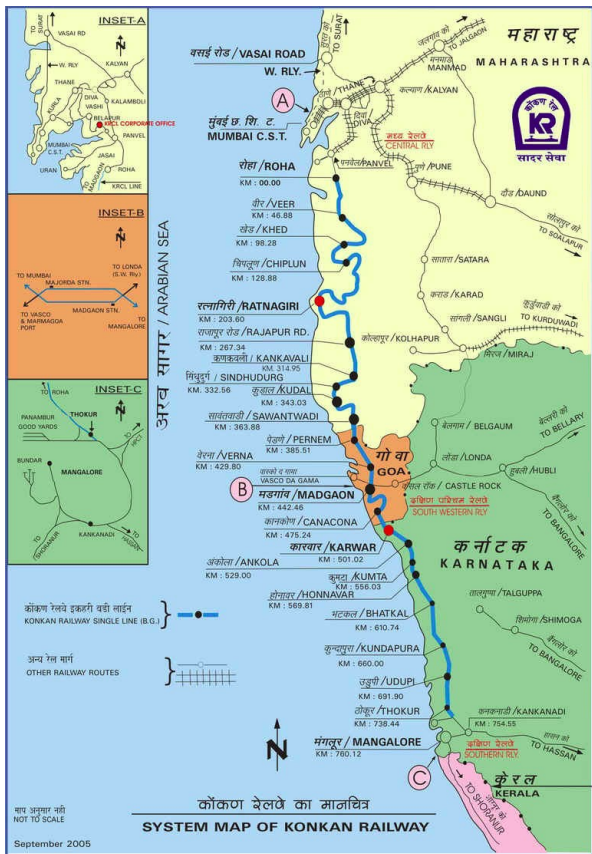
4

Key takeaways from past work

1. Use domain knowledge
 - to divide the problem into a sequence of tasks
 - to define how system performance is measured
2. Define tasks that can be repeatedly performed to achieve goals (constant I/O size)
3. Build the right fidelity of simulation to compute the effect of actions on the system
4. Use RL only for decisions where the 'correct' ones are not obvious
5. Wherever feasible, speed up RL training by seeding with existing heuristics

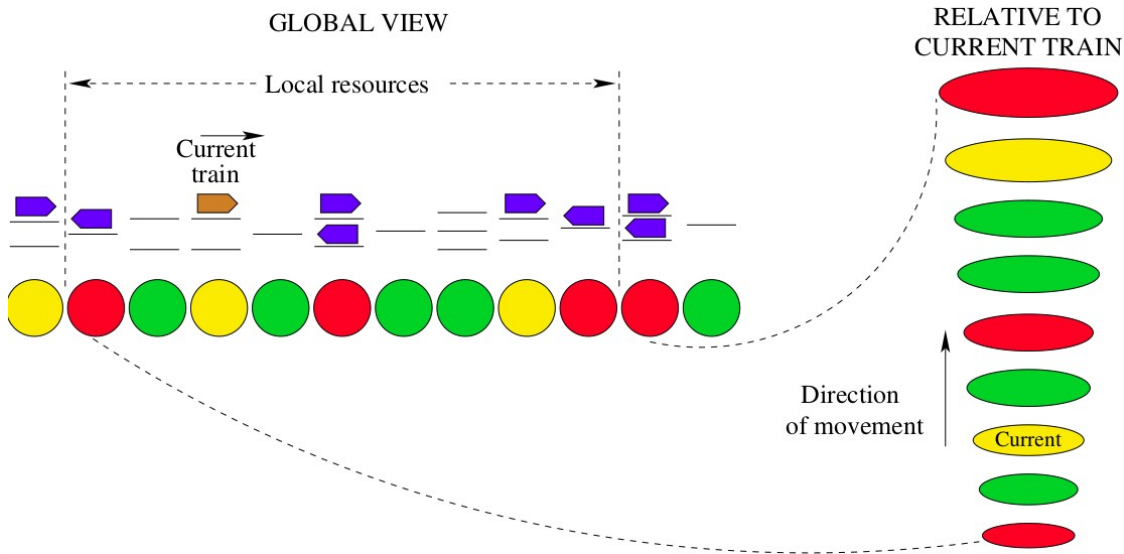
4 RL in the real world

Concrete example: Railway scheduling



Goal: Minimise knock-on effects along the railway line, when recovering from a delayed state

Solution: Divide the problem into a sequence of moves





Current work

5 Current work

Planning for robotic parcel loading

ONLINE 3D BIN PACKING

Stream of incoming boxes



Stable arrangement

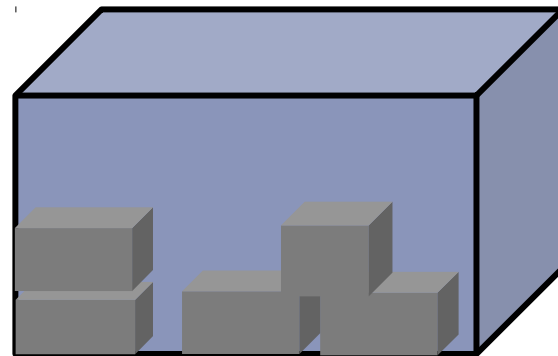
Robot stackable

Goal: Maximise the volume packed in containers, using boxes appearing on a conveyor belt

Rotate the box?

Where to place?

Skip current box?



Container

5 Current work

Supply chain replenishment

OPTIMAL NETWORK OPERATION



Goal: Minimise supply chain operating costs while maximising key performance indicators

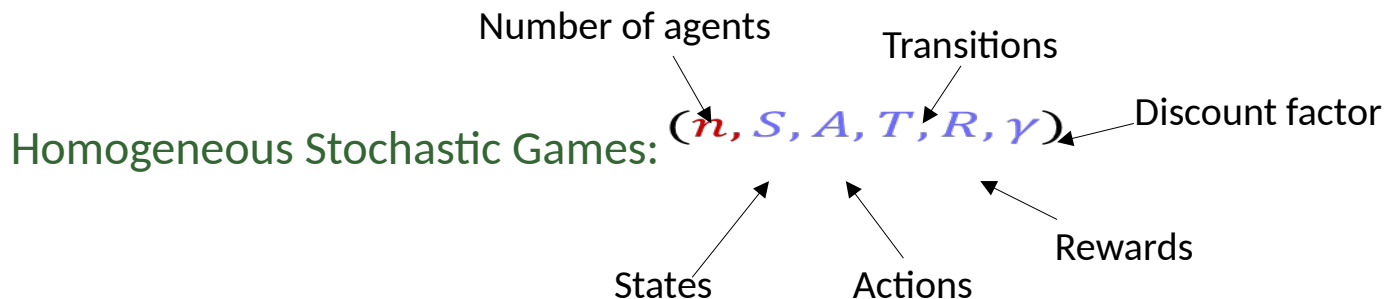
State of the art: Requirements flow from right to left (upstream), while products flow from left to right (downstream)

Solution: Multi-agent reinforcement learning at each node of the supply chain, for automated adaptive response to demands

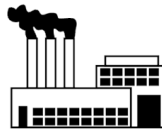
Avoid stock-out
Reduce wastage

Replenishment decisions
Transport & labour planning

Generalisation of Markov Decision Processes to **Stochastic Games**



Heterogeneous Stochastic Games:



Can this set of participants in a **system of systems** collaborate effectively?

Reinforcement learning = Use of machine learning for decision-making problems

Should be used when it is the best tool for the job:

1. Fast response
2. Systems *simulatable* but not analytically describable
3. Unknown 'optimal' decisions
4. Sequence-dependent rewards

Making RL work for you in real life:

1. Make sure you can simulate your problem, for training
2. Divide large problems into a sequence of repeated tasks
 3. Use domain expertise rather than throw it away
 4. Build solutions with explanations, not black boxes