

On optimal dynamic treatment regimes (full reinforcement learning)

Nilanjana Laha



Broad goal

Month 1



Month 2



Month 3



These images are generated by Dall-E

Chronic diseases demand ongoing treatments. Can we apply reinforcement learning for optimal, **patient-specific**, data-driven treatment policy?

Where does it stand as an area?

Where does it stand as an area?

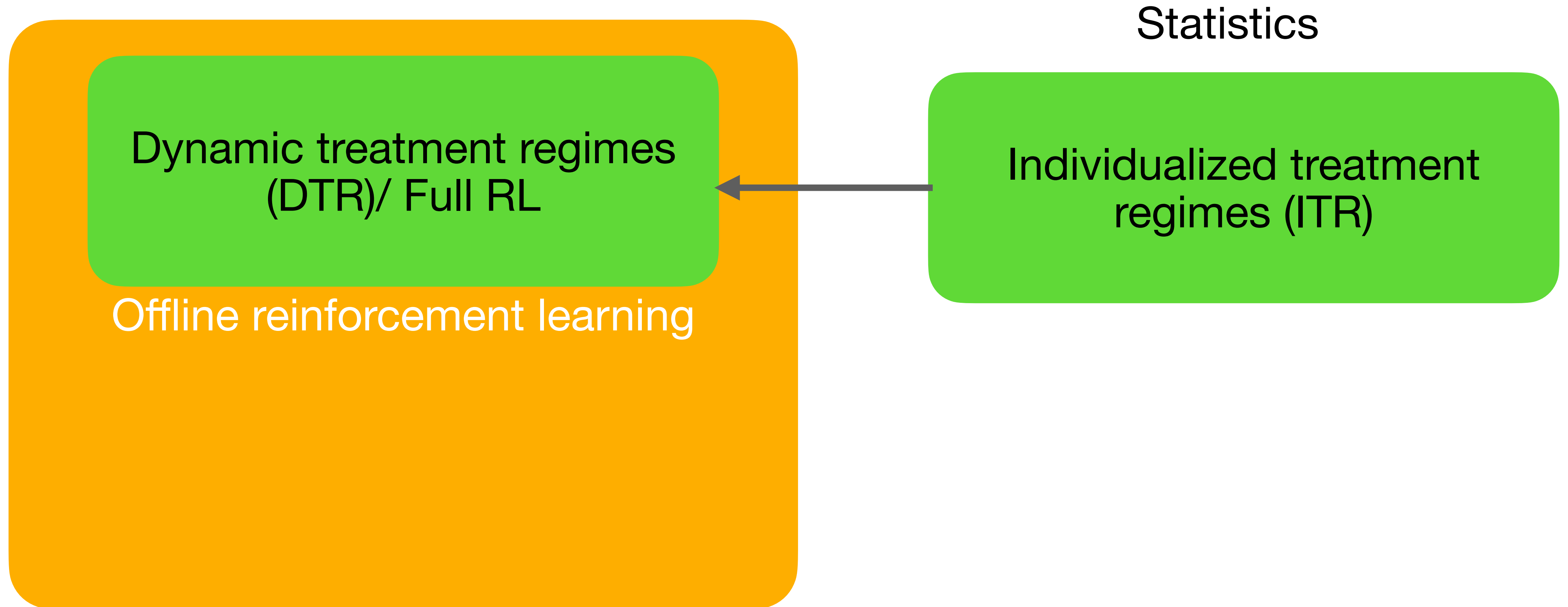
Dynamic treatment regimes
(DTR)/ Full RL

Where does it stand as an area?

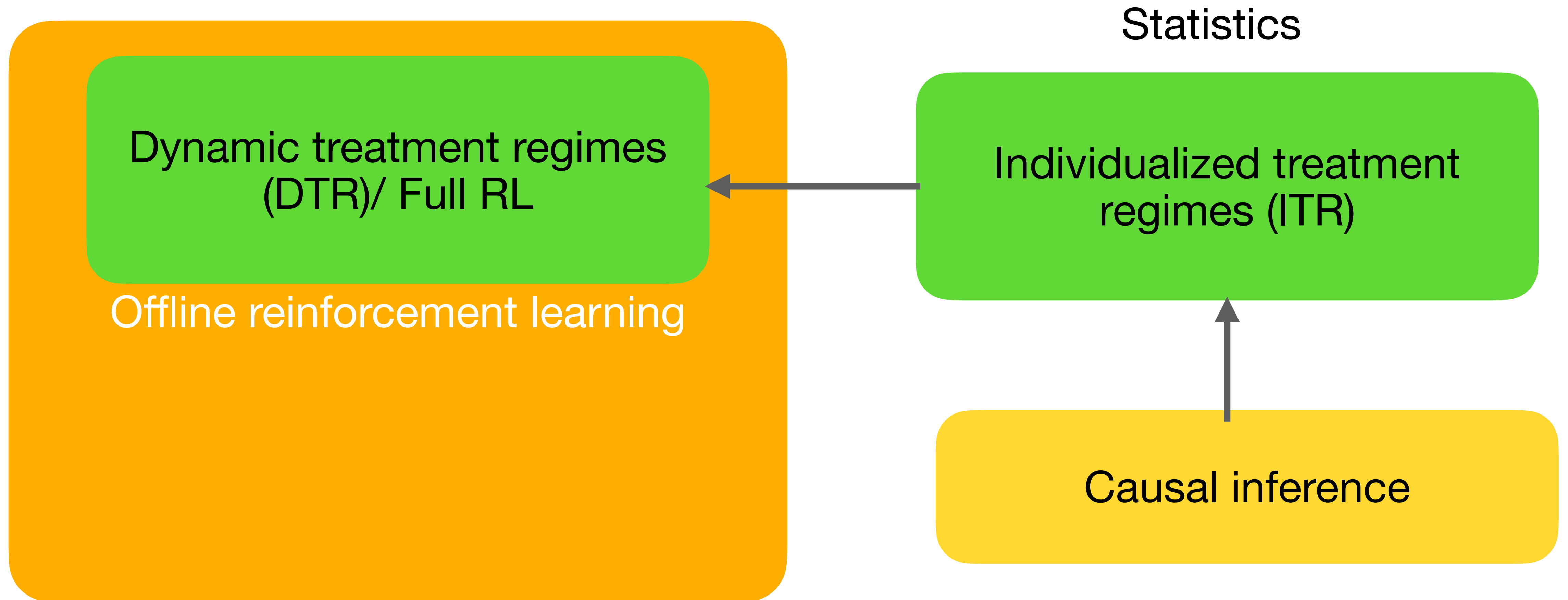
Dynamic treatment regimes
(DTR)/ Full RL

Offline reinforcement learning

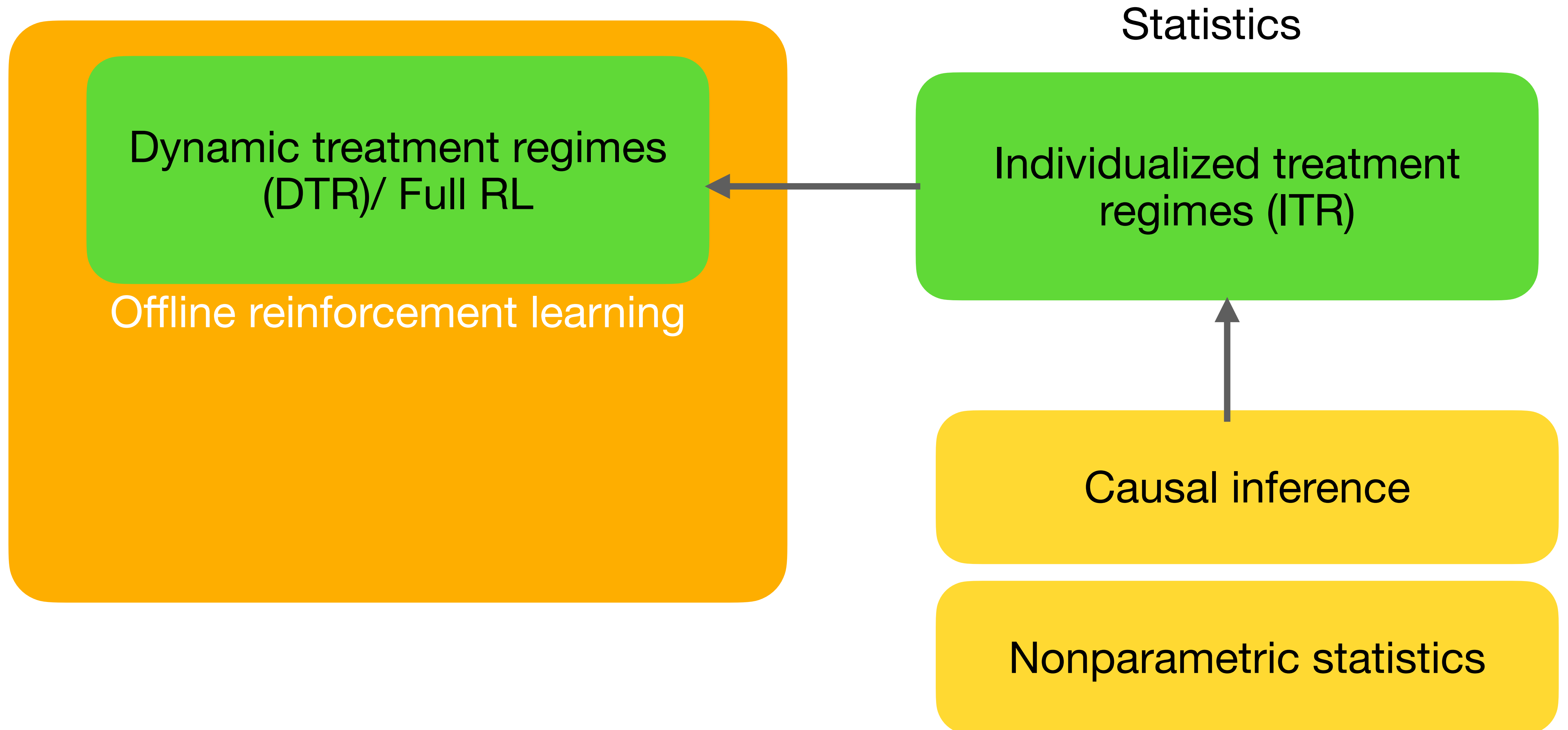
Where does it stand as an area?



Where does it stand as an area?



Where does it stand as an area?



Outline

- Example: sepsis
- Problem formulation
- Proposed method
- Open questions

Example: sepsis

Sepsis

Sepsis

Cause: Body's response to infection injures own tissues, organs.



Image source: MedicineNet

Sepsis

Cause: Body's response to infection injures own tissues, organs.



Image source: MedicineNet

Expensive

In-patient cost > \$22 billion

Sepsis

Cause: Body's response to infection injures own tissues, organs.



Image source: MedicineNet

Expensive

In-patient cost > \$22 billion

Challenging

Fatality 30 %

Sepsis

Cause: Body's response to infection injures own tissues, organs.



Image source: MedicineNet

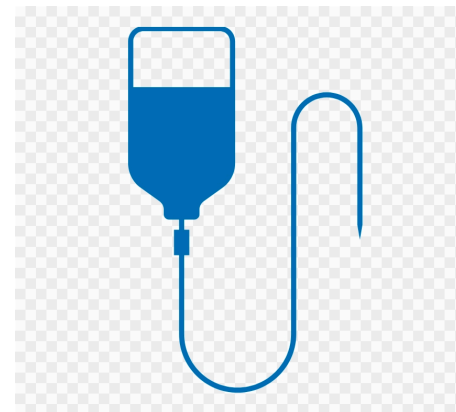
Expensive

In-patient cost > \$22 billion

Challenging

Fatality 30 %

Popular treatment:



**IV-fluid
administration**

Sepsis

Cause: Body's response to infection injures own tissues, organs.



Image source: MedicineNet

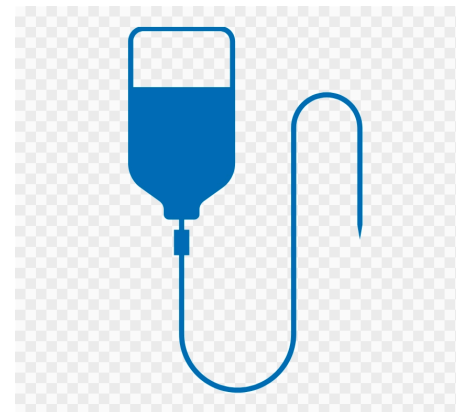
Expensive

In-patient cost > \$22 billion

Challenging

Fatality 30 %

Popular treatment:



**IV-fluid
administration**

Goal: policy learning for IV-fluid administration

Sepsis-3 data (Beth Israel Hospital, Boston)

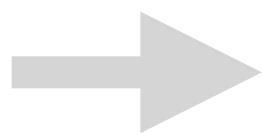
Sepsis-3 data (Beth Israel Hospital, Boston)

Hour 0

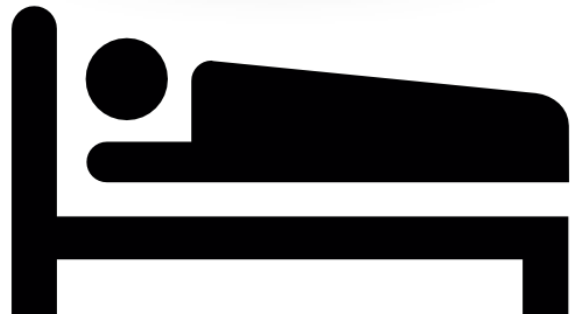


Sepsis-3 data (Beth Israel Hospital, Boston)

Hour 0



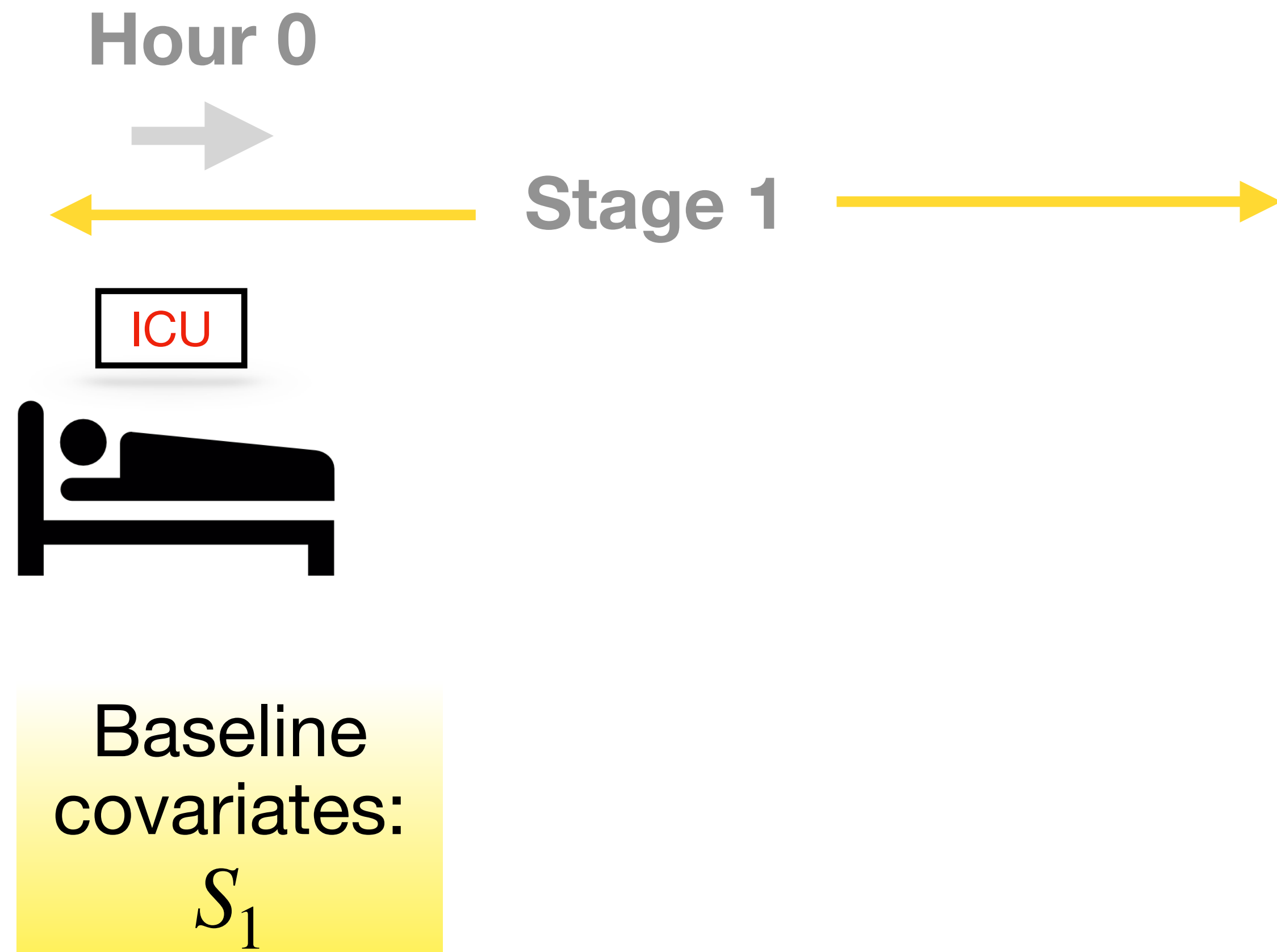
ICU



Sepsis-3 data (Beth Israel Hospital, Boston)



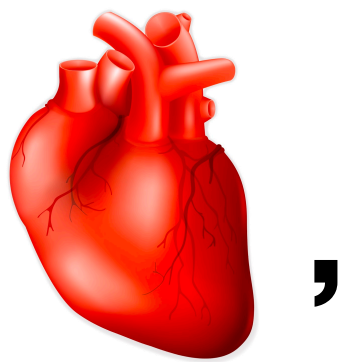
Sepsis-3 data (Beth Israel Hospital, Boston)



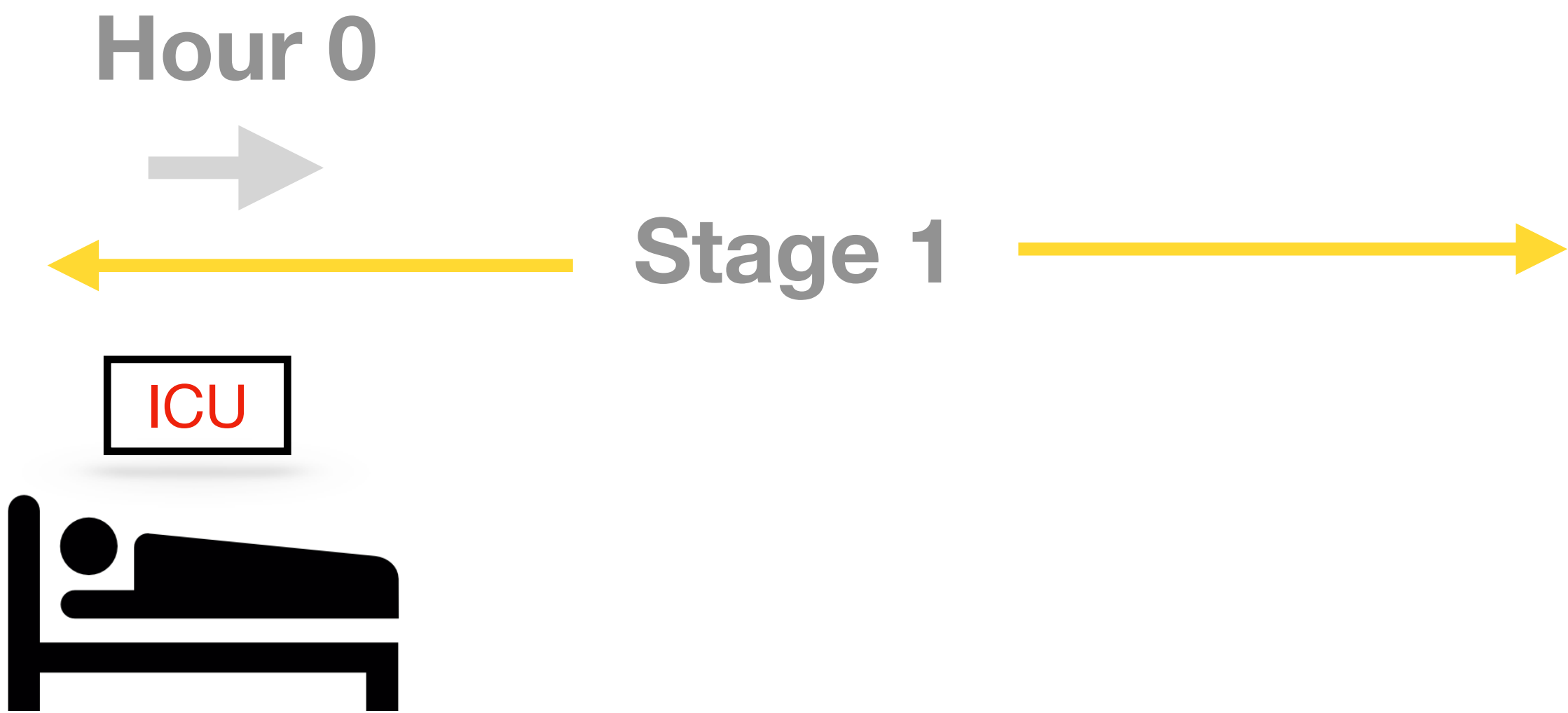
Sepsis-3 data (Beth Israel Hospital, Boston)



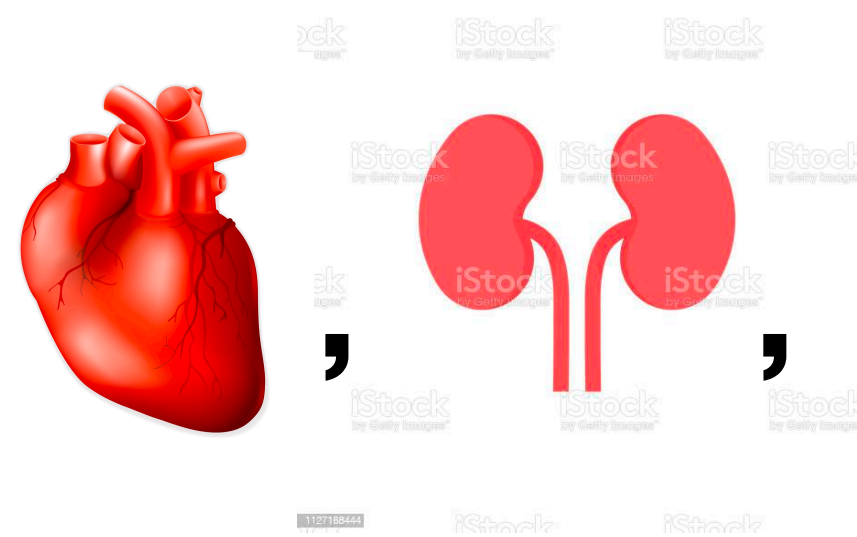
Baseline
covariates:
 S_1



Sepsis-3 data (Beth Israel Hospital, Boston)



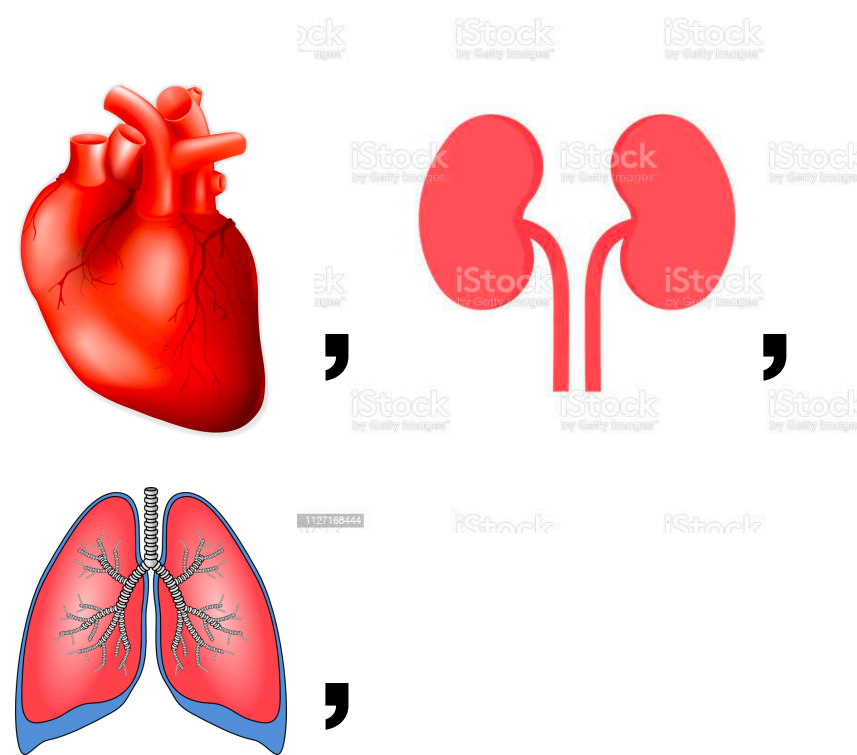
Baseline
covariates:
 S_1



Sepsis-3 data (Beth Israel Hospital, Boston)



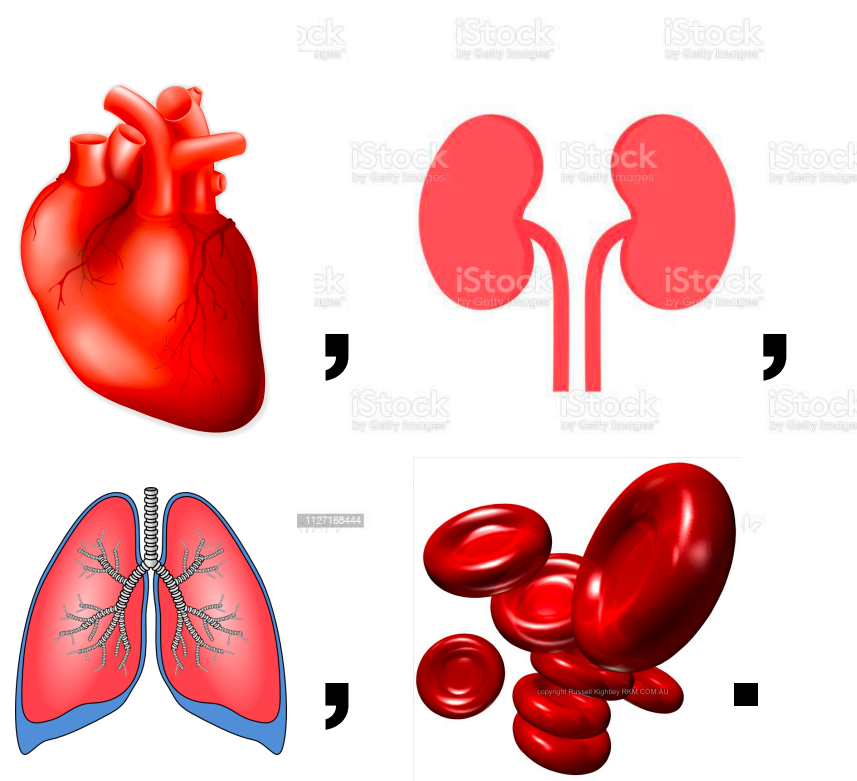
Baseline
covariates:
 S_1



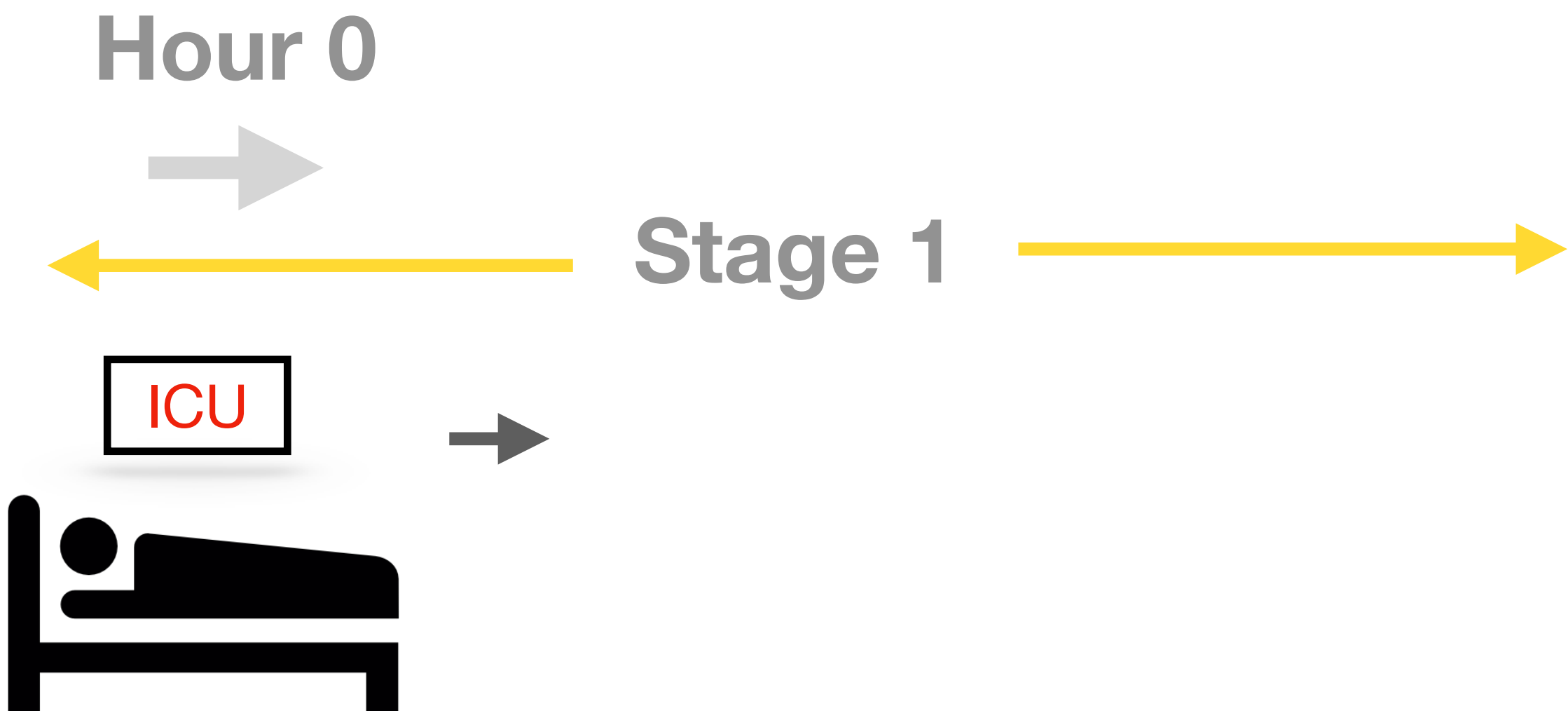
Sepsis-3 data (Beth Israel Hospital, Boston)



Baseline
covariates:
 S_1

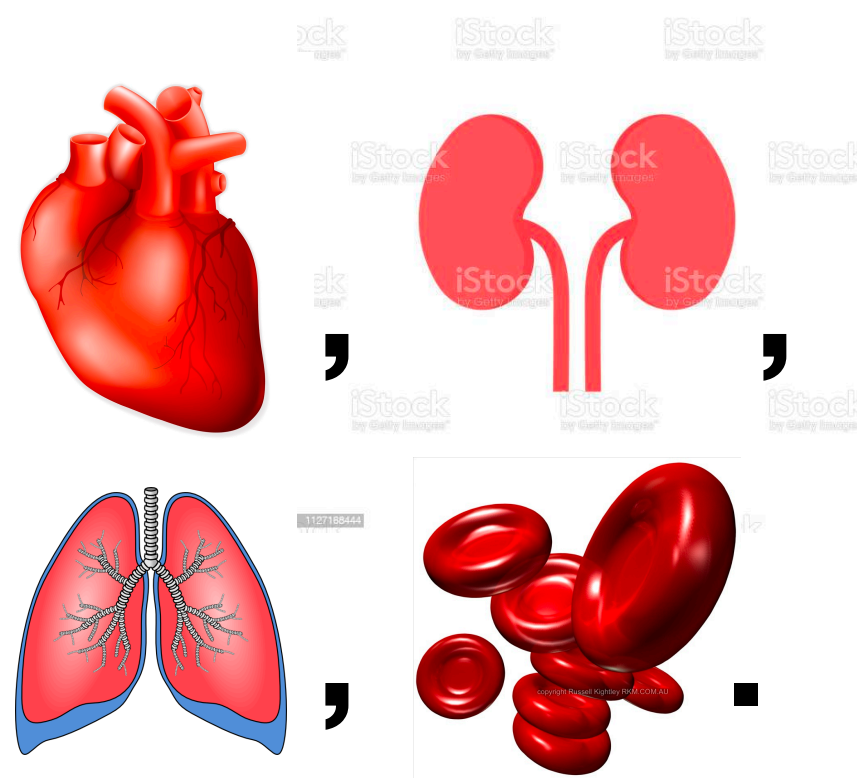


Sepsis-3 data (Beth Israel Hospital, Boston)

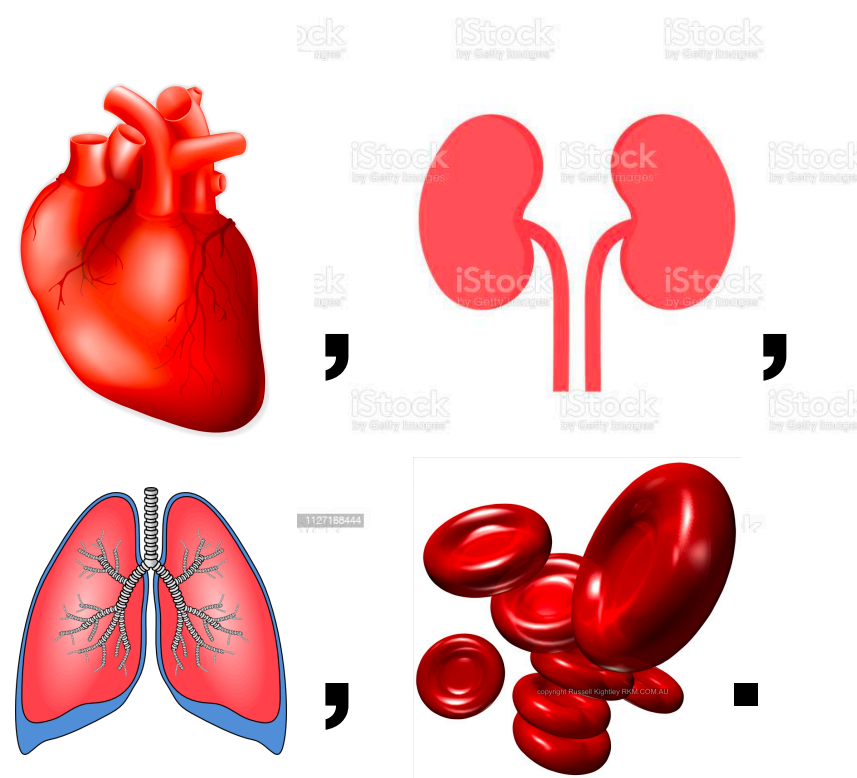
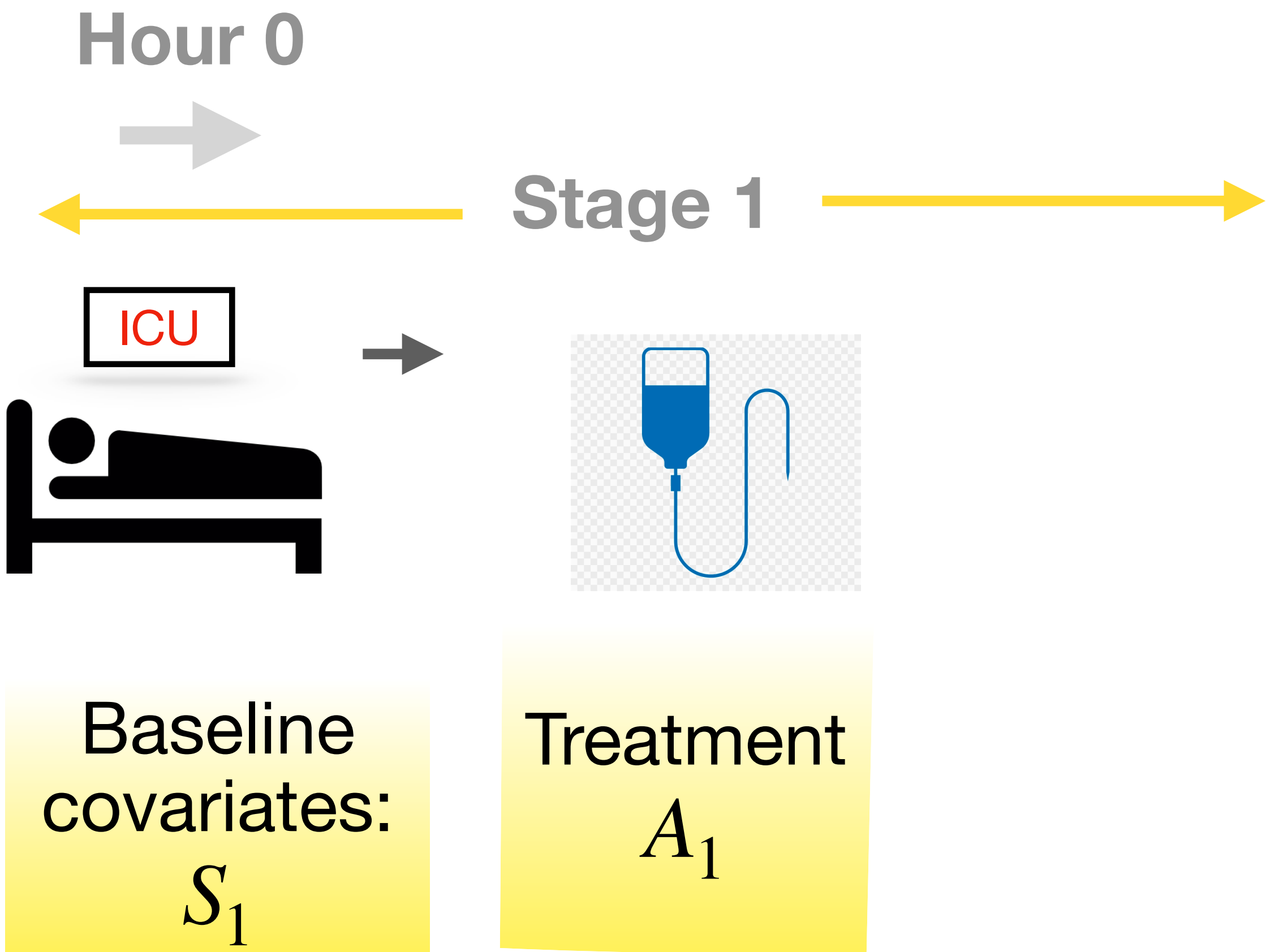


Baseline
covariates:
 S_1

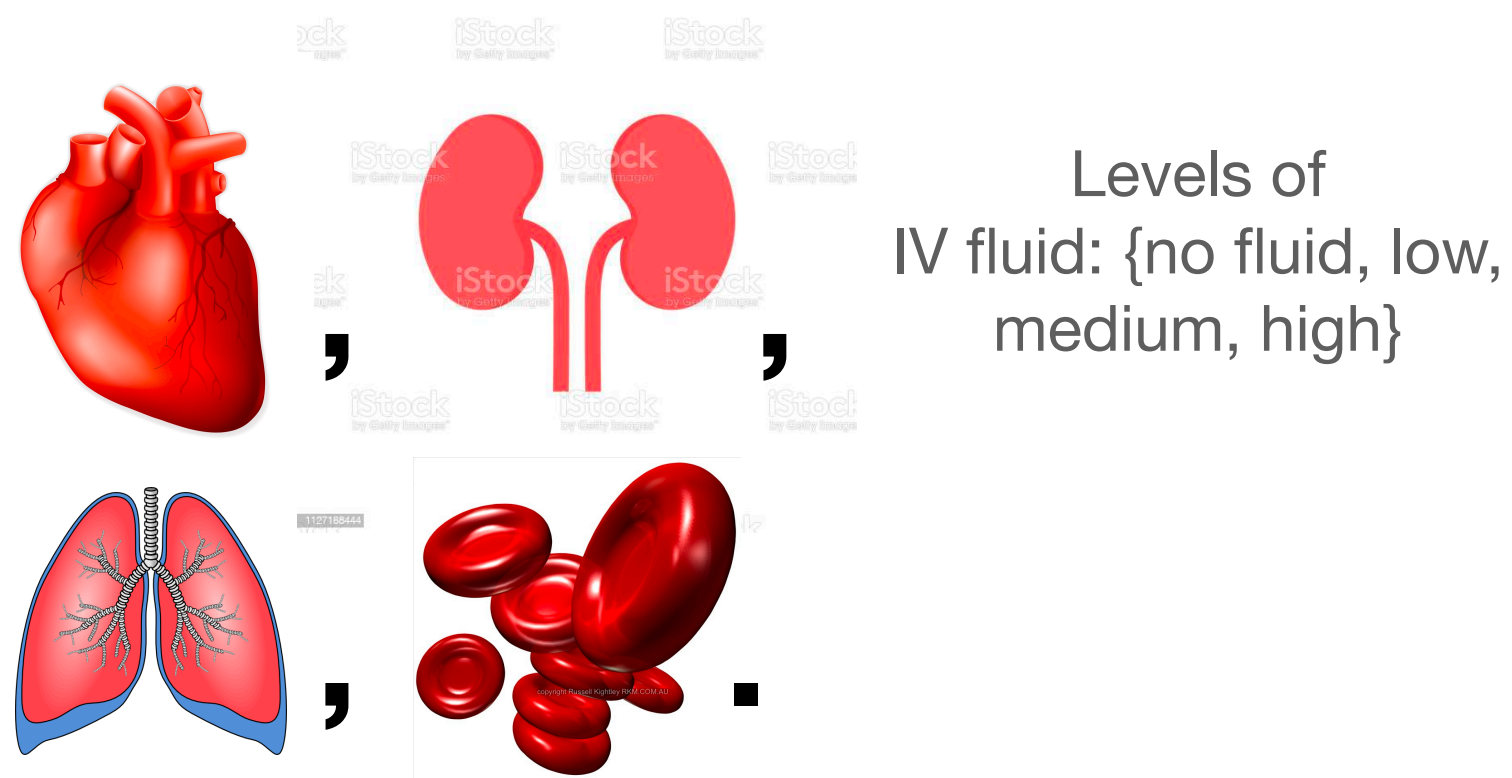
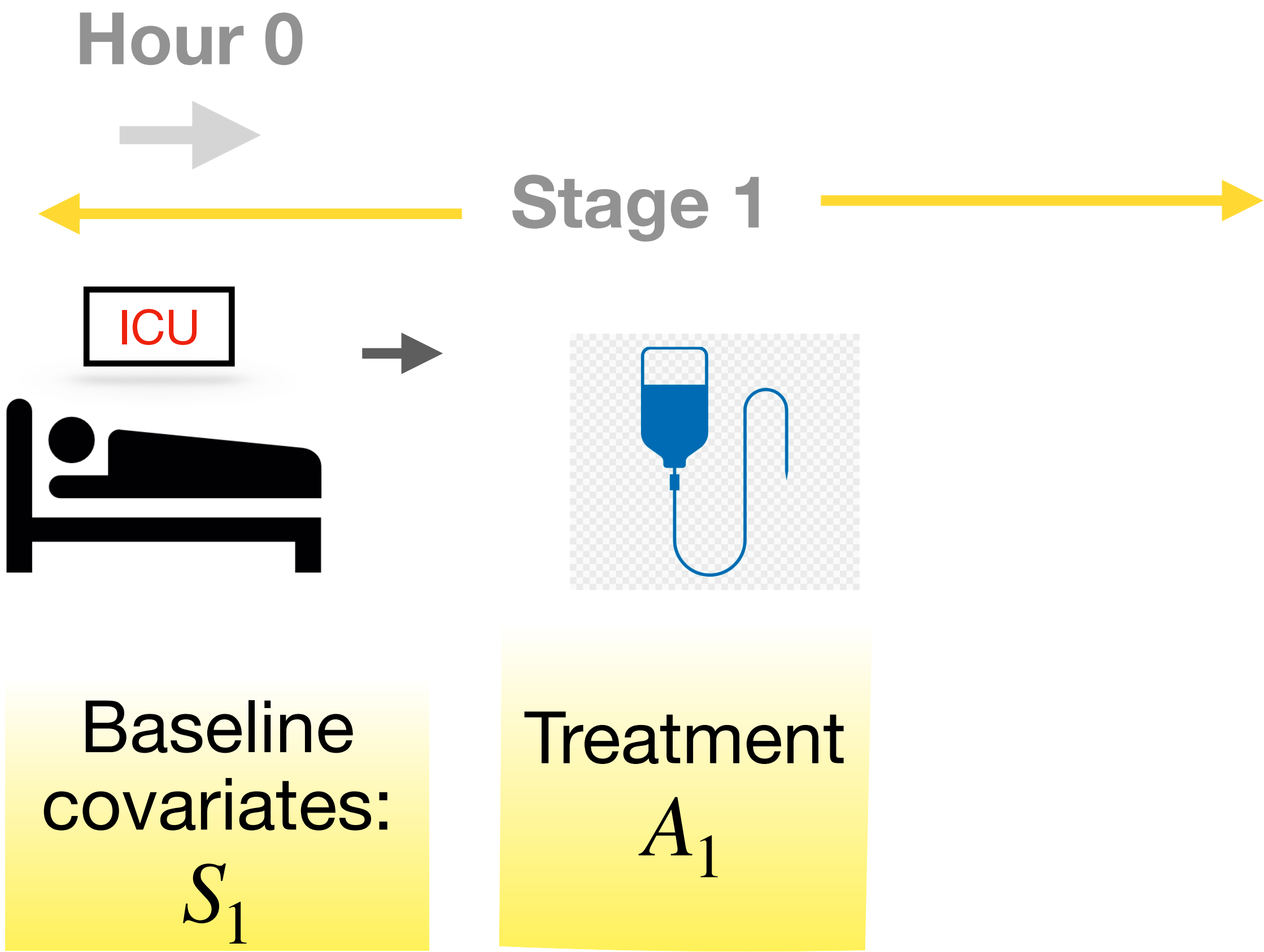
Treatment
 A_1



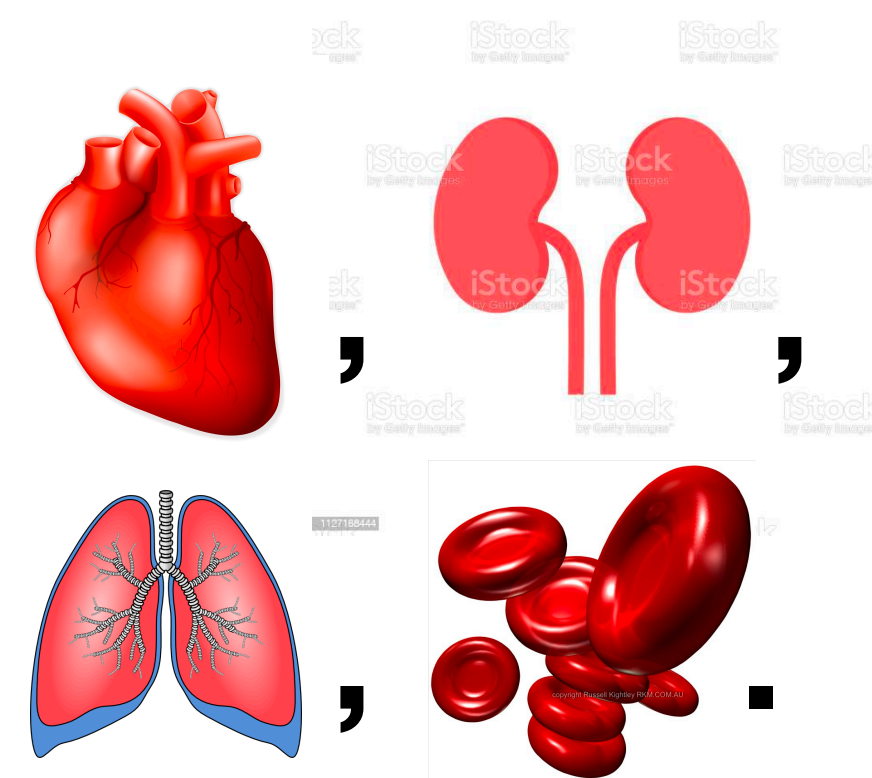
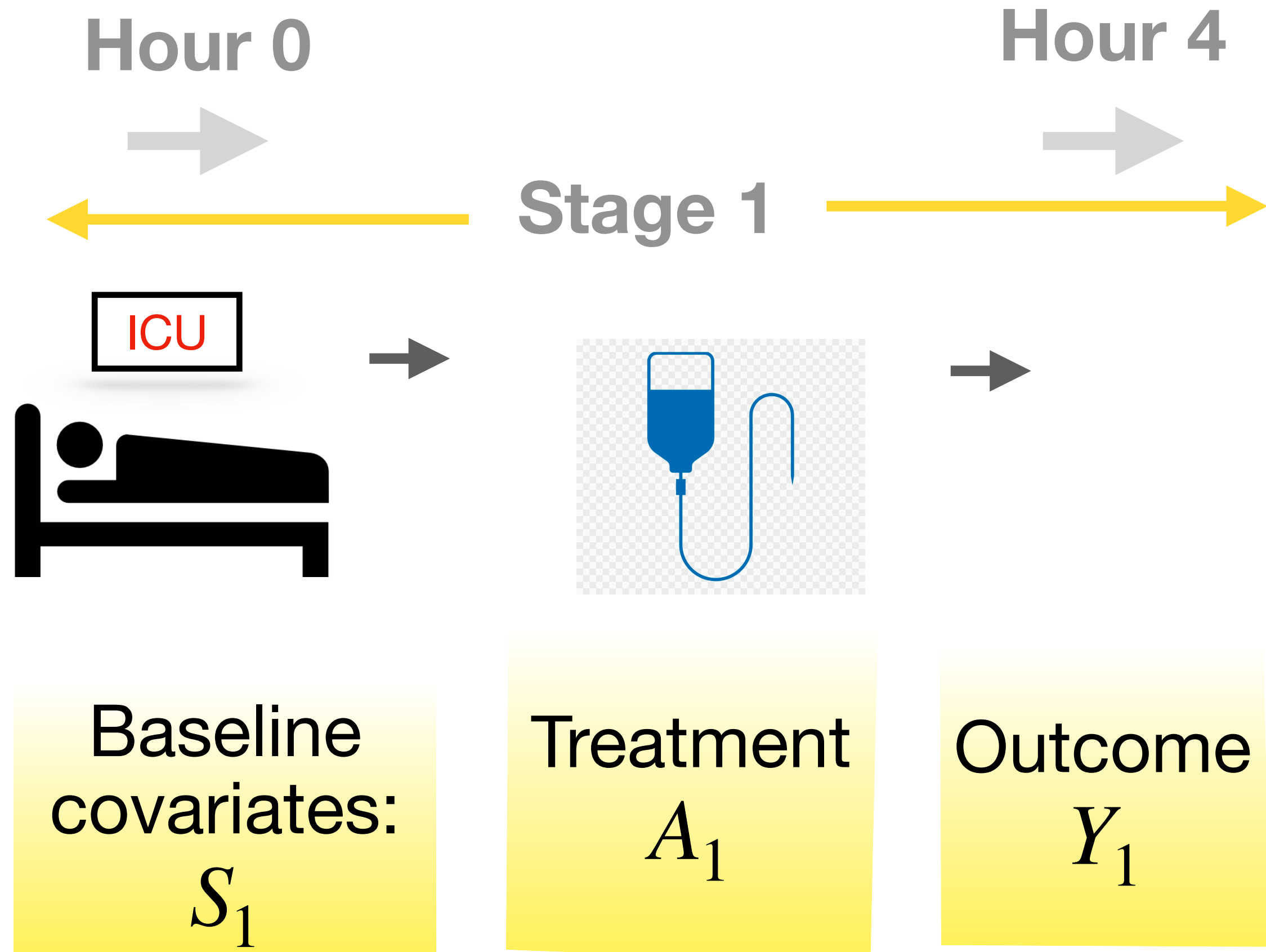
Sepsis-3 data (Beth Israel Hospital, Boston)



Sepsis-3 data (Beth Israel Hospital, Boston)

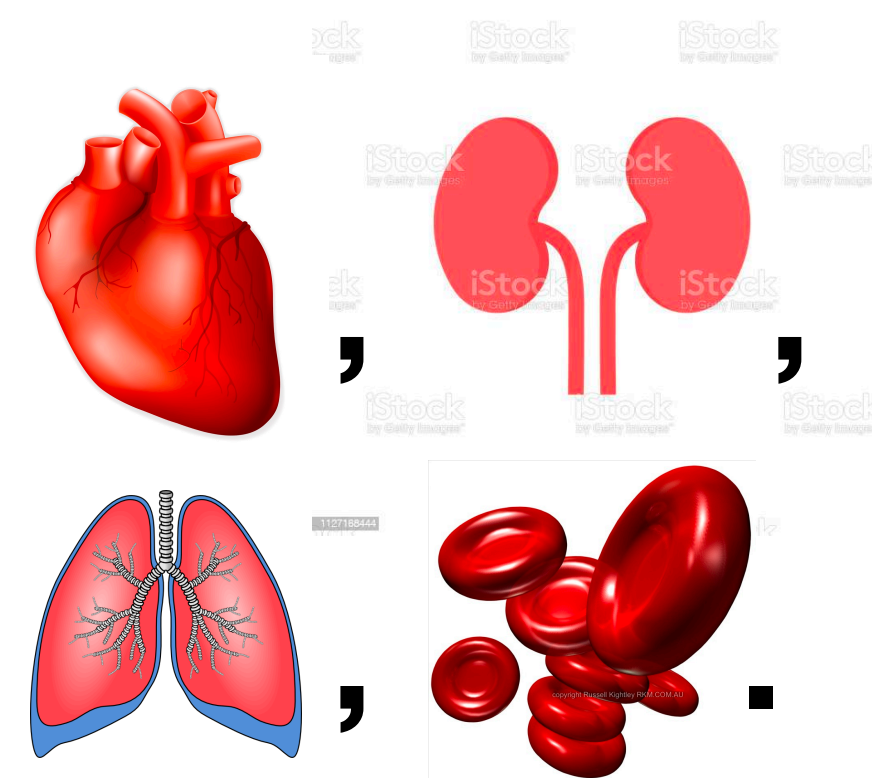
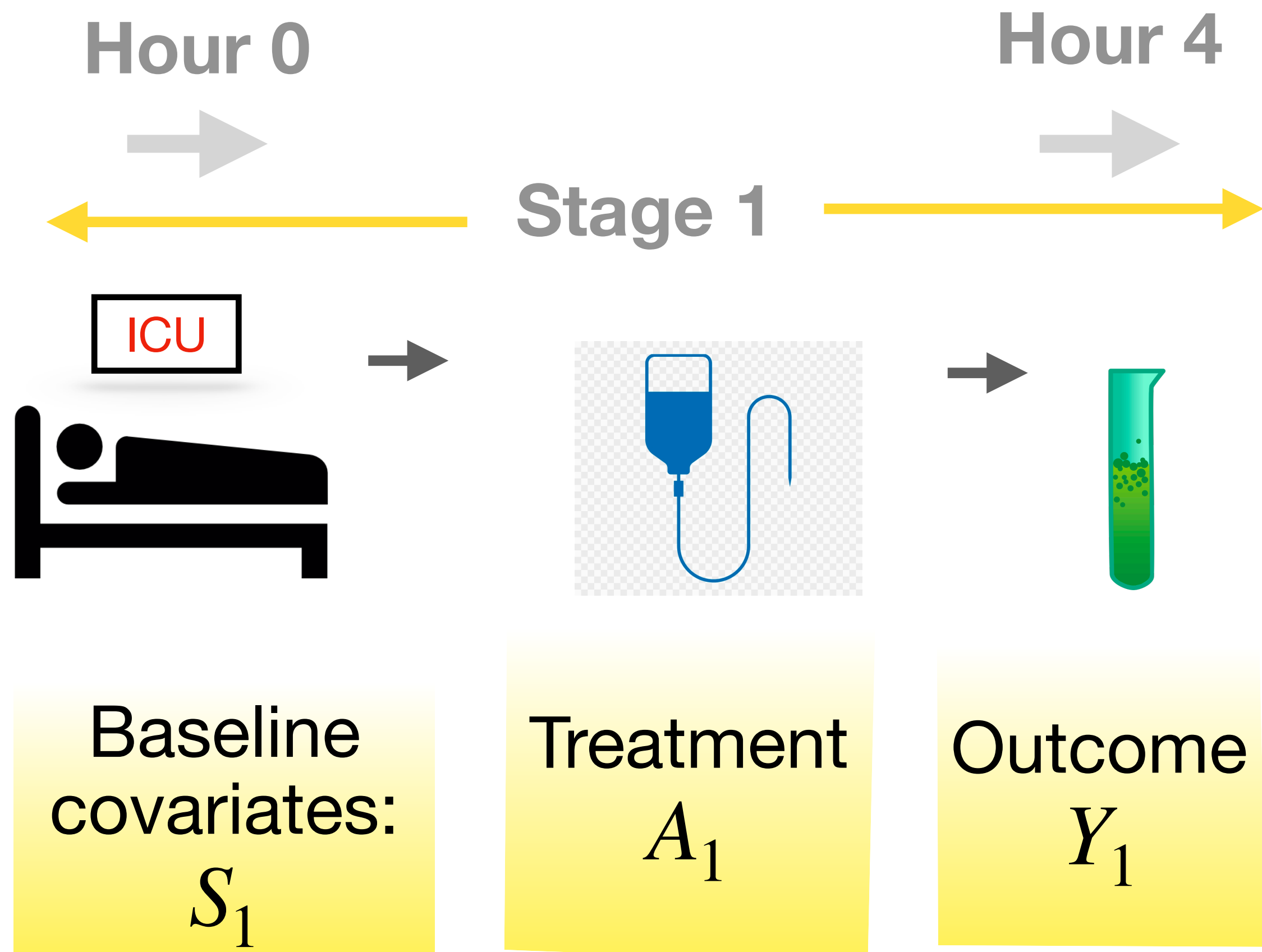


Sepsis-3 data (Beth Israel Hospital, Boston)



Levels of
IV fluid: {no fluid, low,
medium, high}

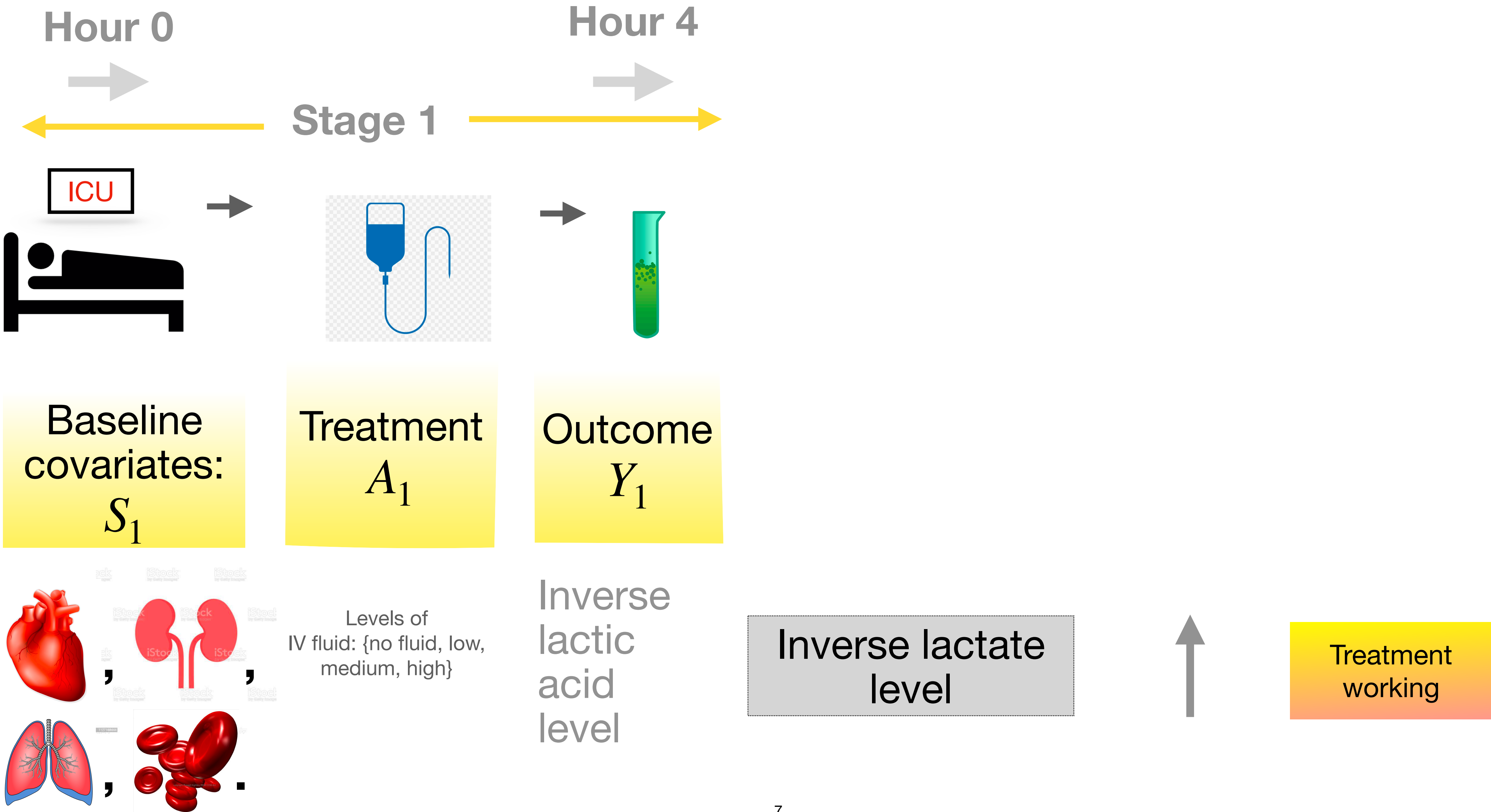
Sepsis-3 data (Beth Israel Hospital, Boston)



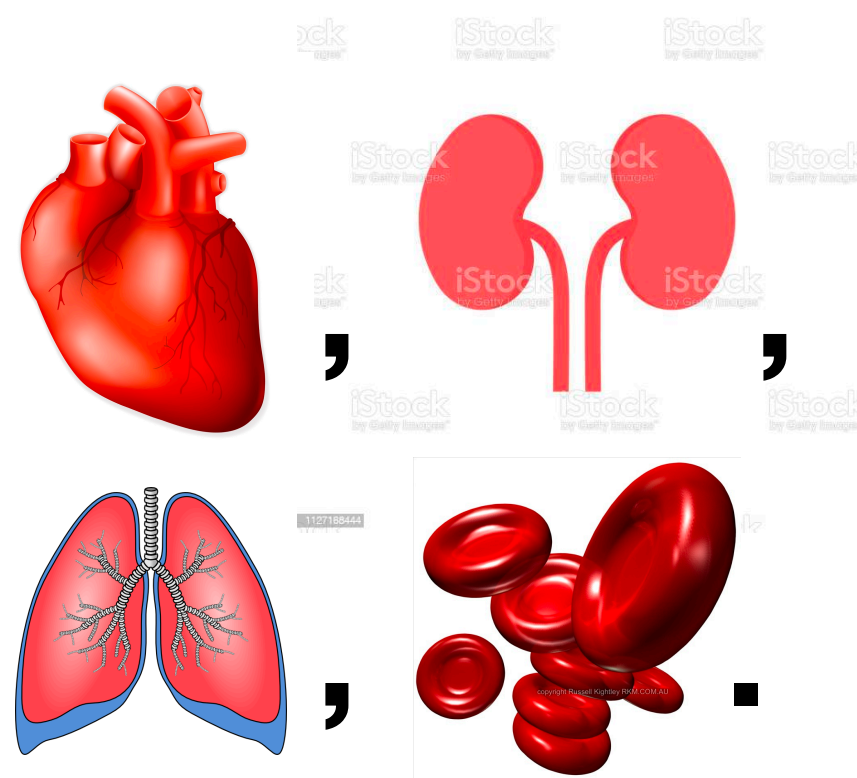
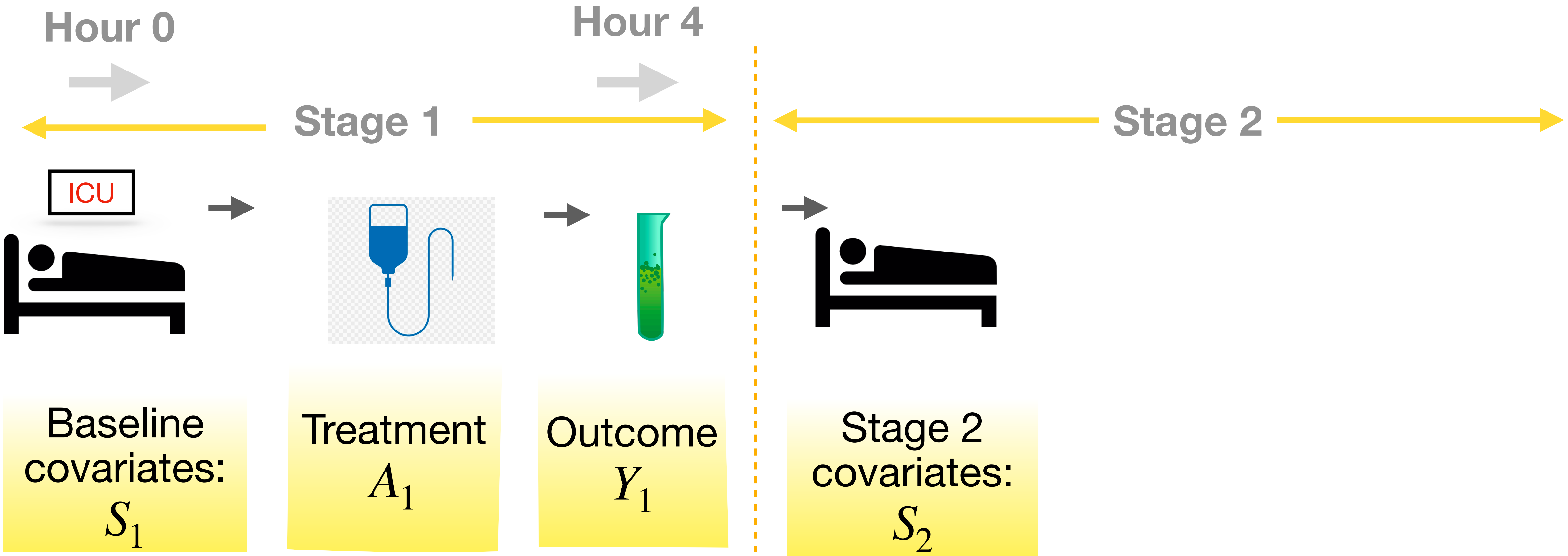
Levels of
IV fluid: {no fluid, low,
medium, high}

Inverse lactic acid level

Sepsis-3 data (Beth Israel Hospital, Boston)



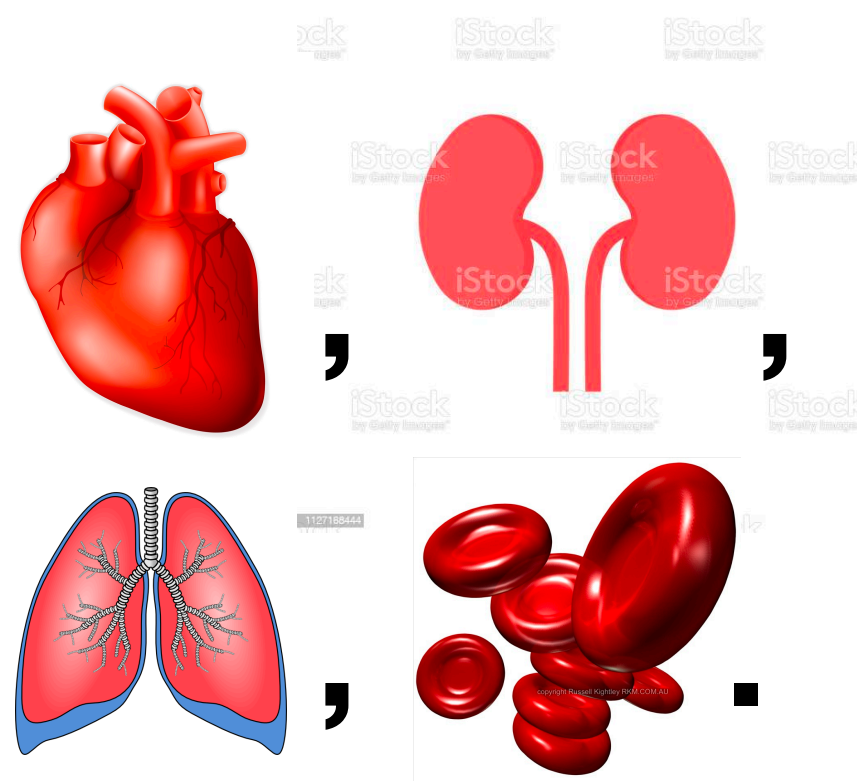
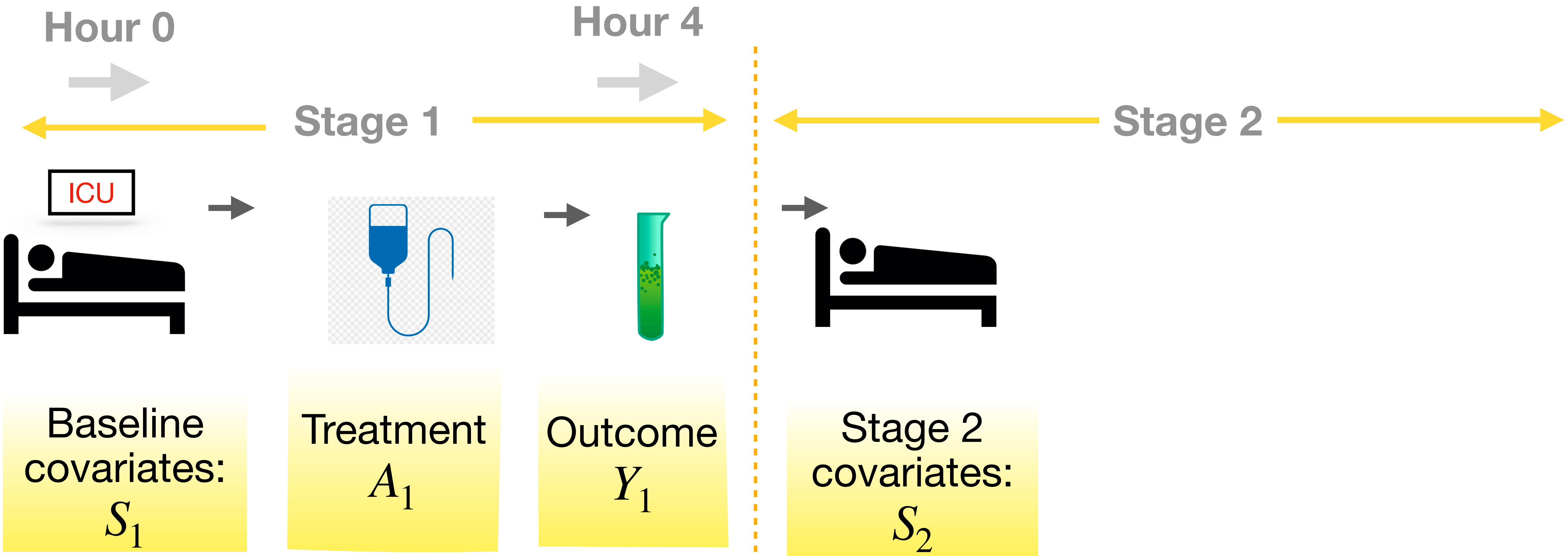
Sepsis-3 data (Beth Israel Hospital, Boston)



Levels of
IV fluid: {no fluid, low,
medium, high}

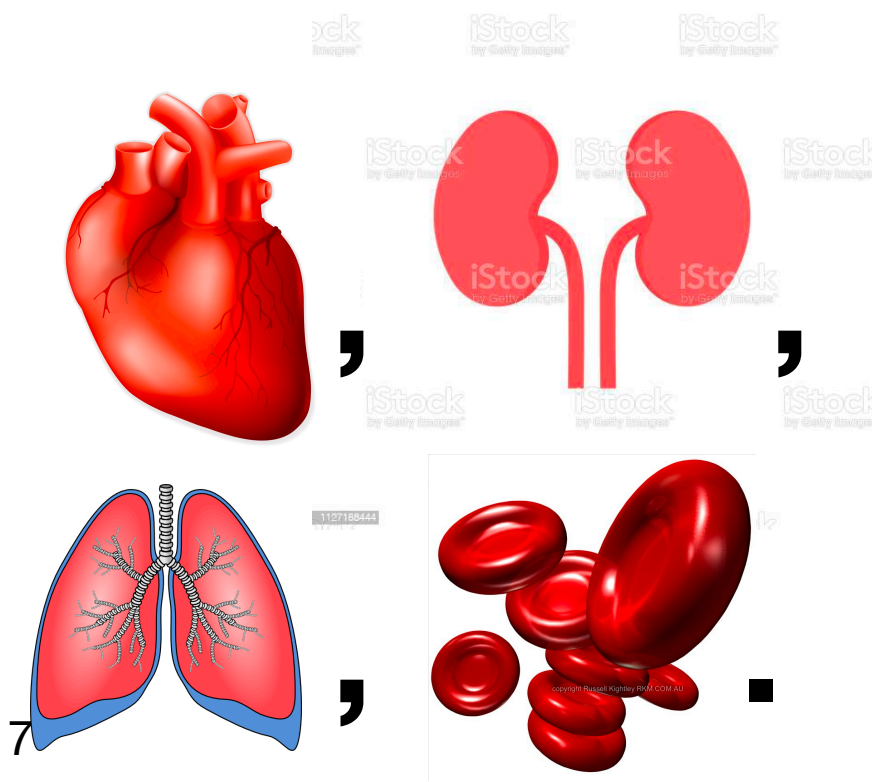
Inverse
lactic
acid
level

Sepsis-3 data (Beth Israel Hospital, Boston)

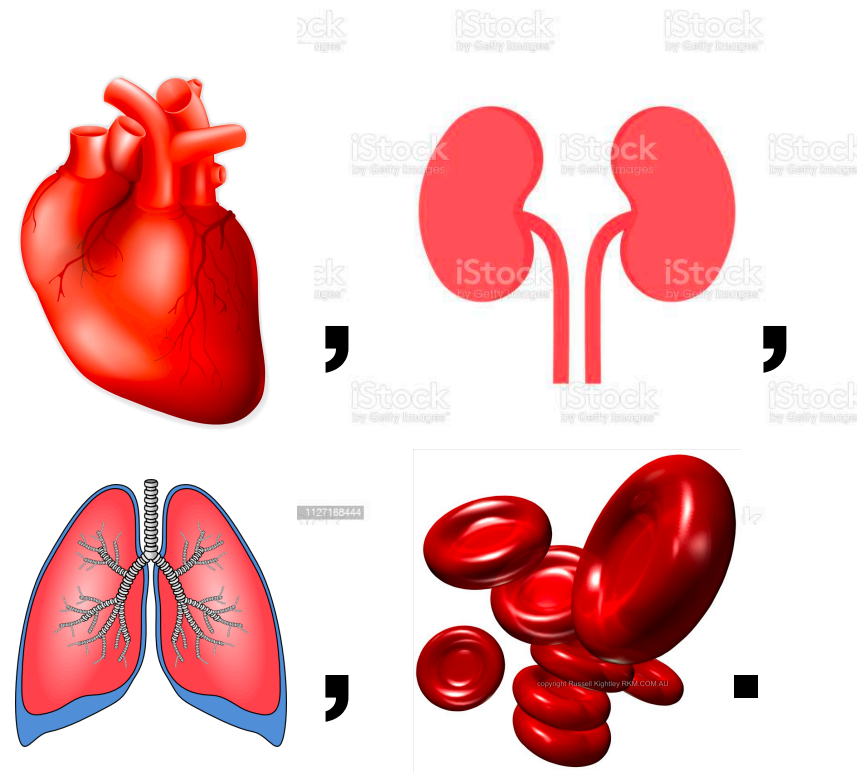
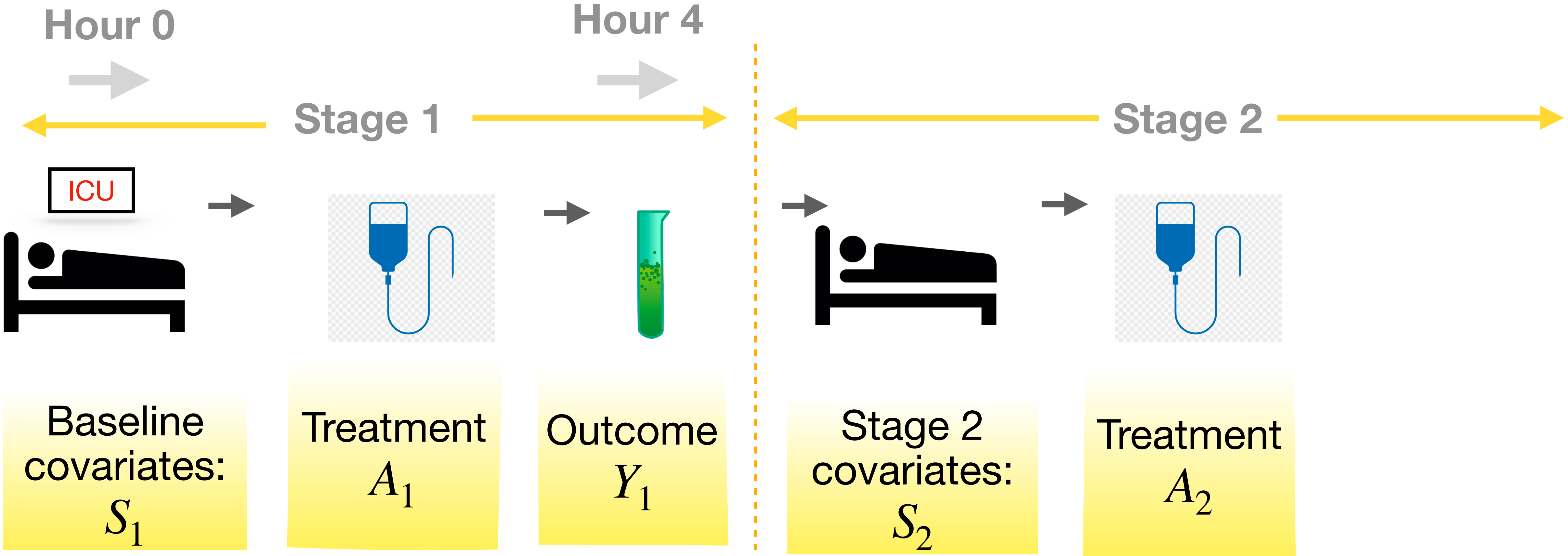


Levels of
IV fluid: {no fluid, low,
medium, high}

Inverse
lactic
acid
level

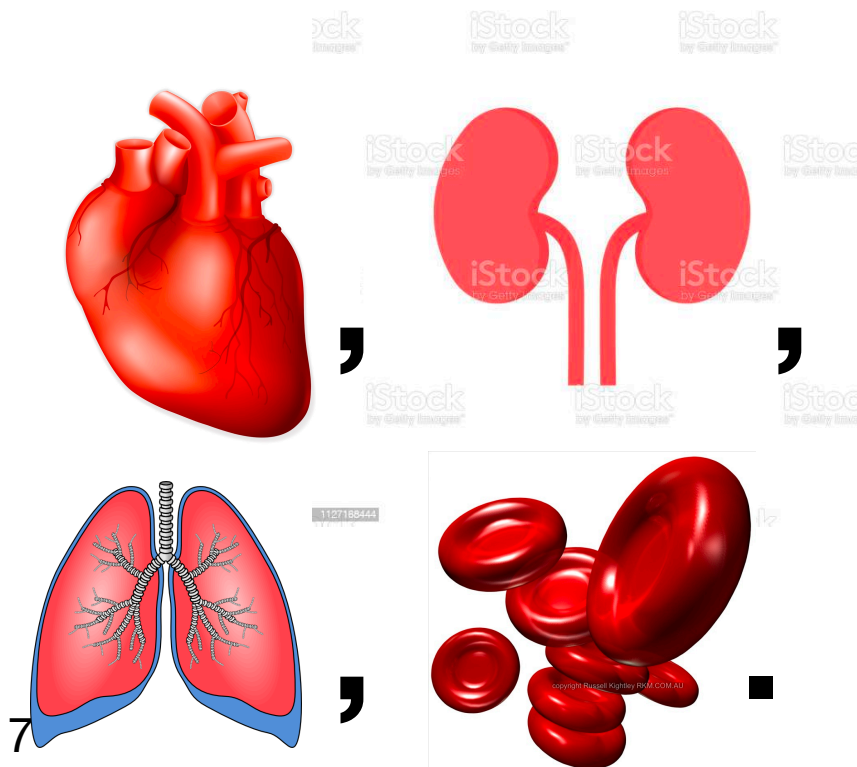


Sepsis-3 data (Beth Israel Hospital, Boston)

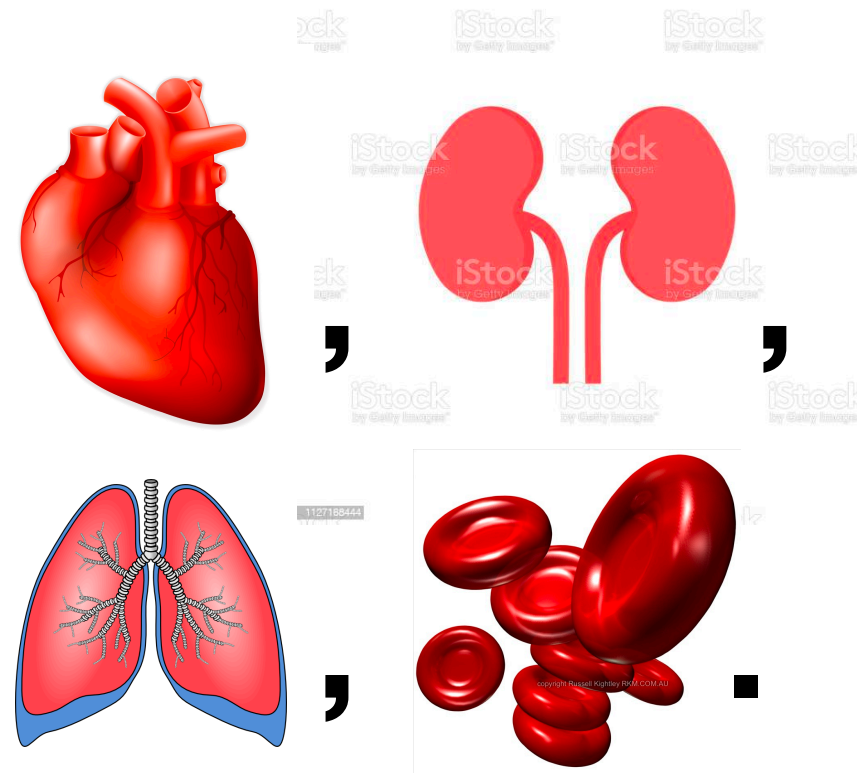
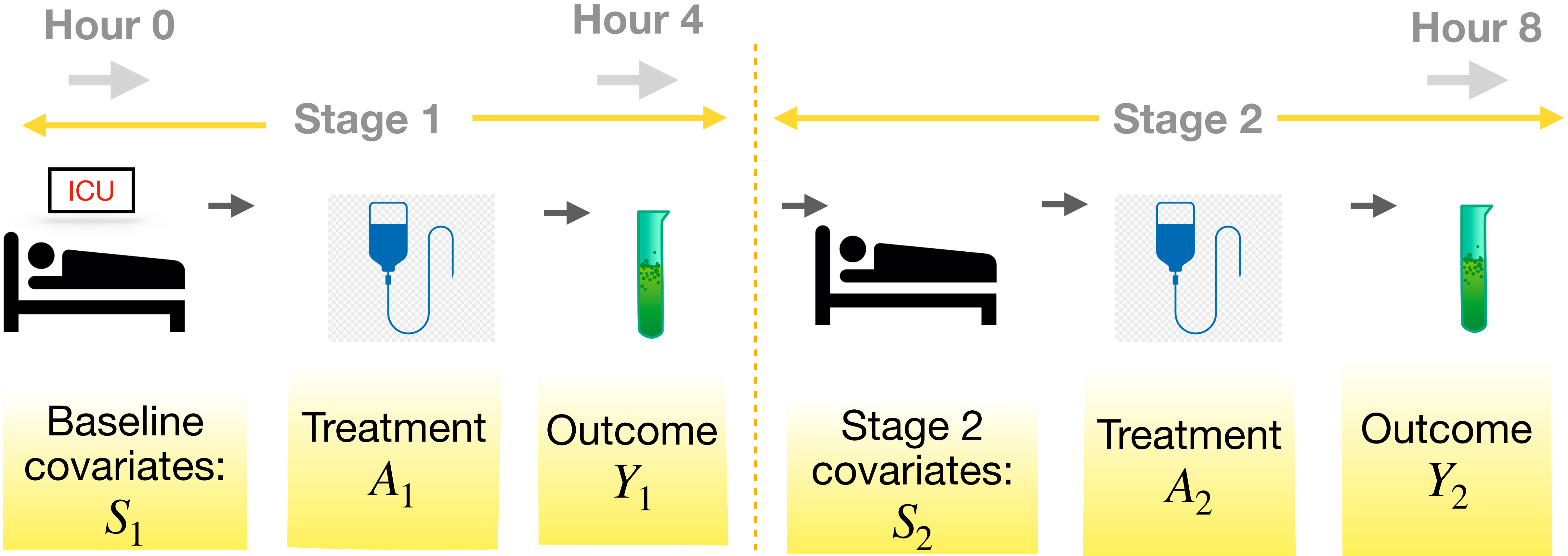


Levels of
IV fluid: {no fluid, low,
medium, high}

Inverse
lactic
acid
level

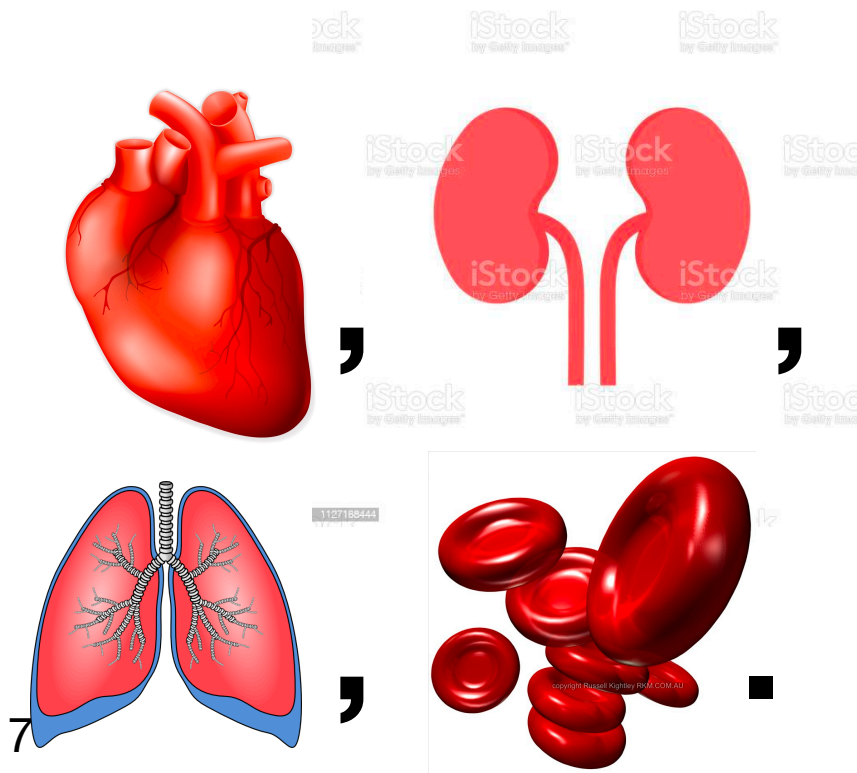


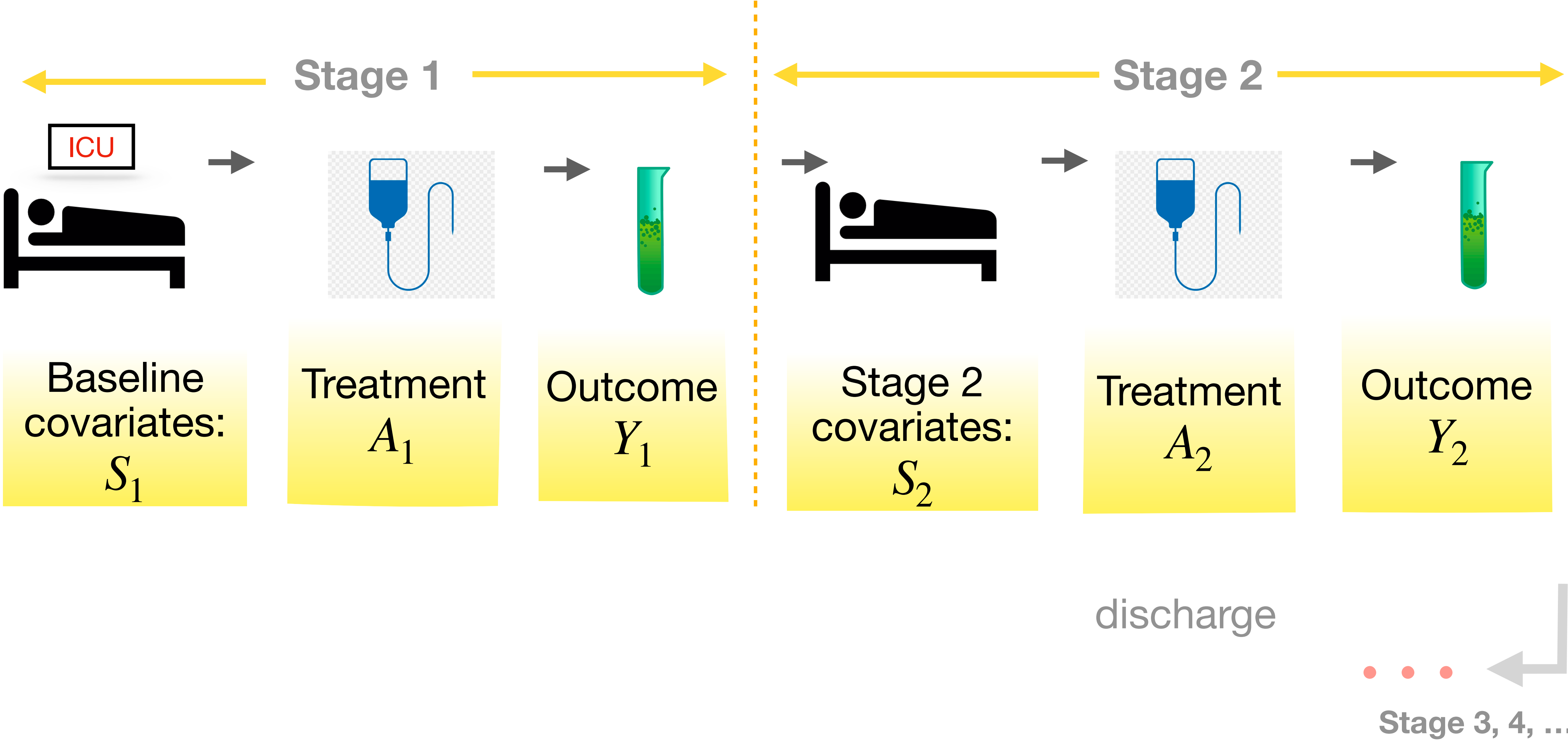
Sepsis-3 data (Beth Israel Hospital, Boston)

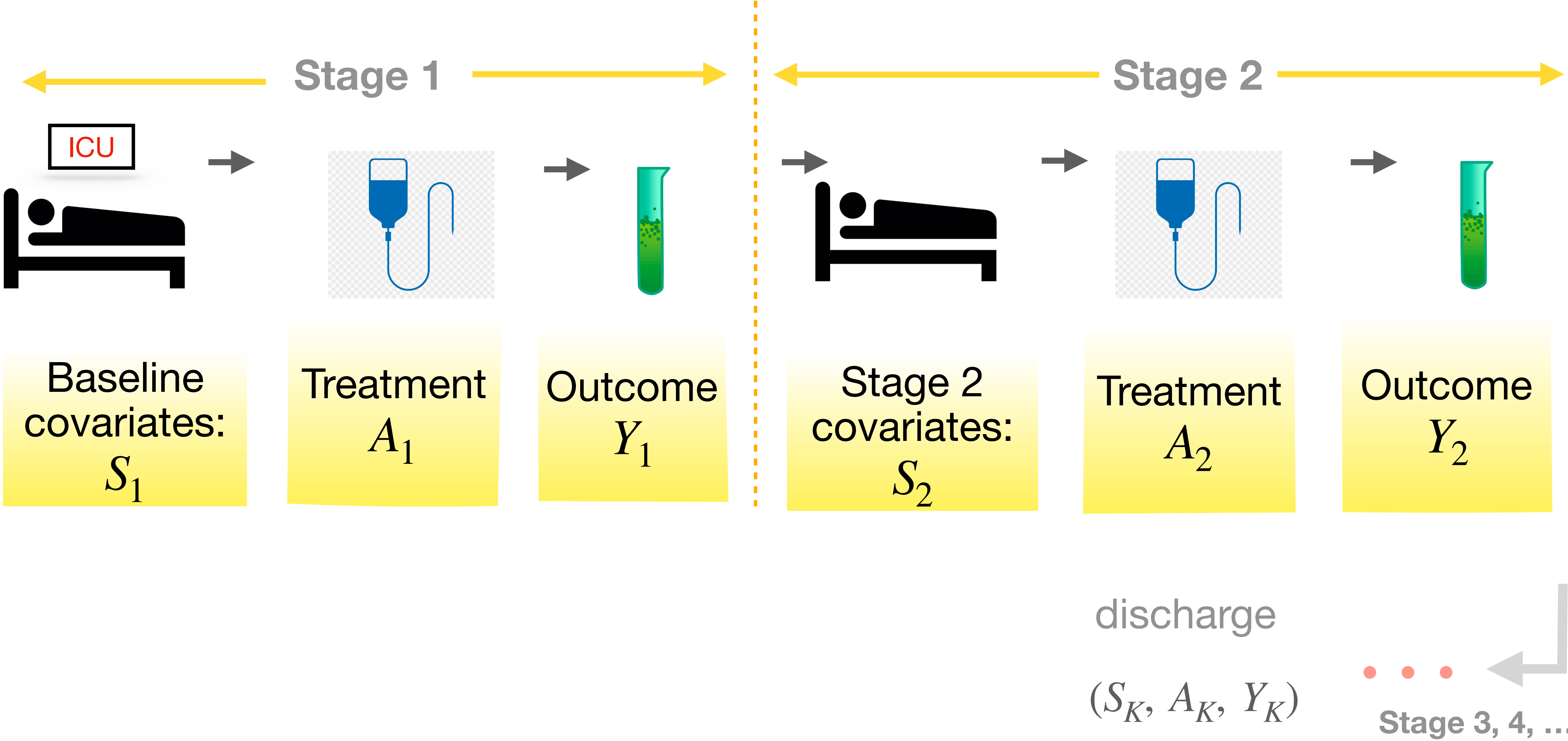


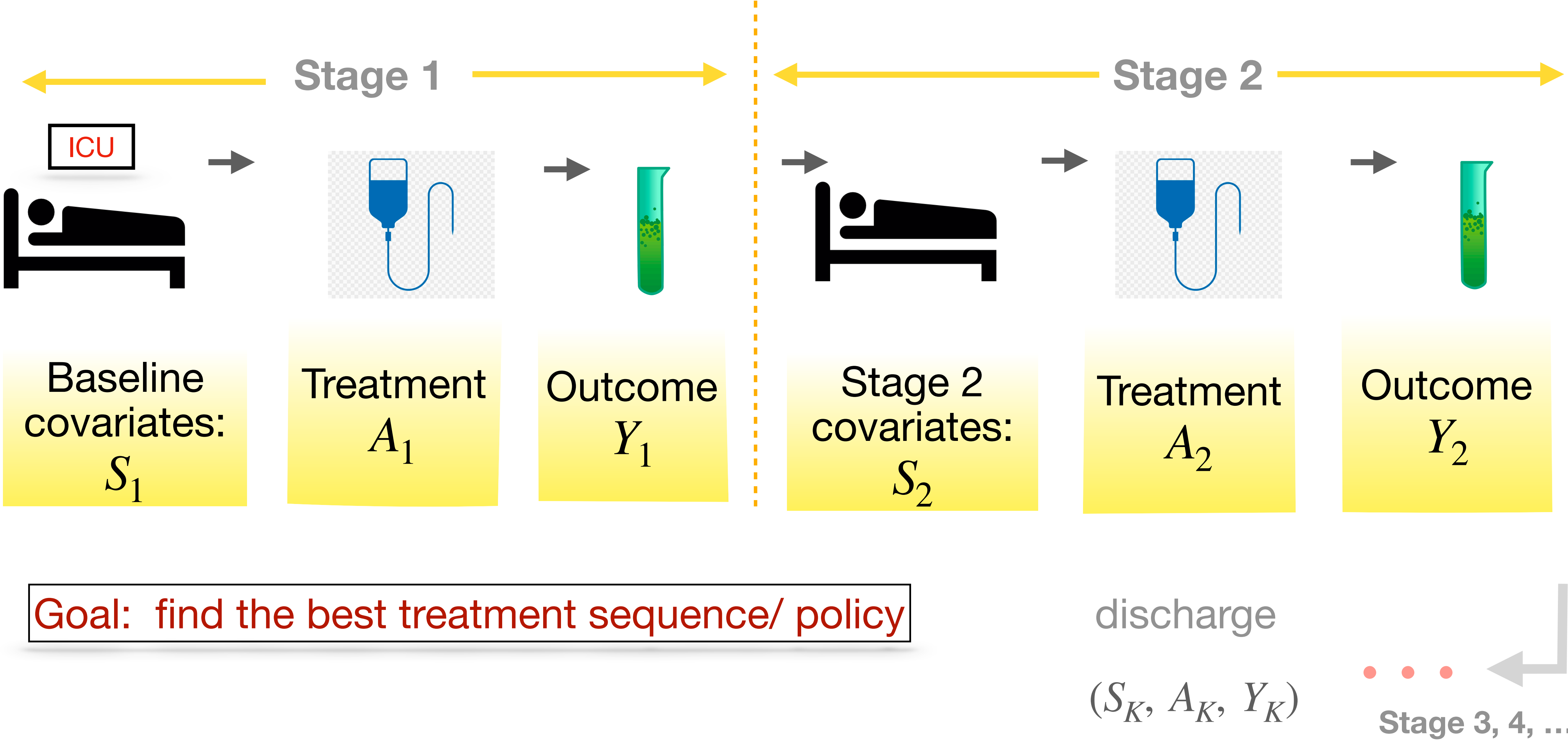
Levels of
IV fluid: {no fluid, low,
medium, high}

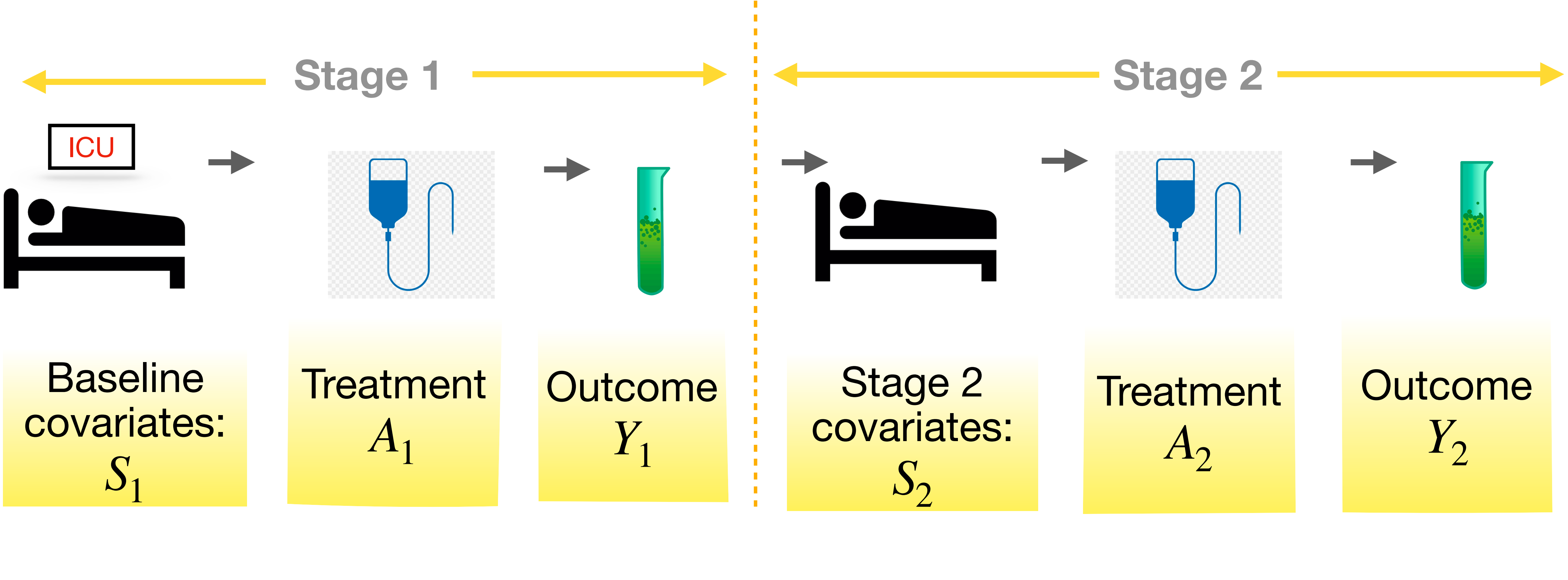
Inverse
lactic
acid
level







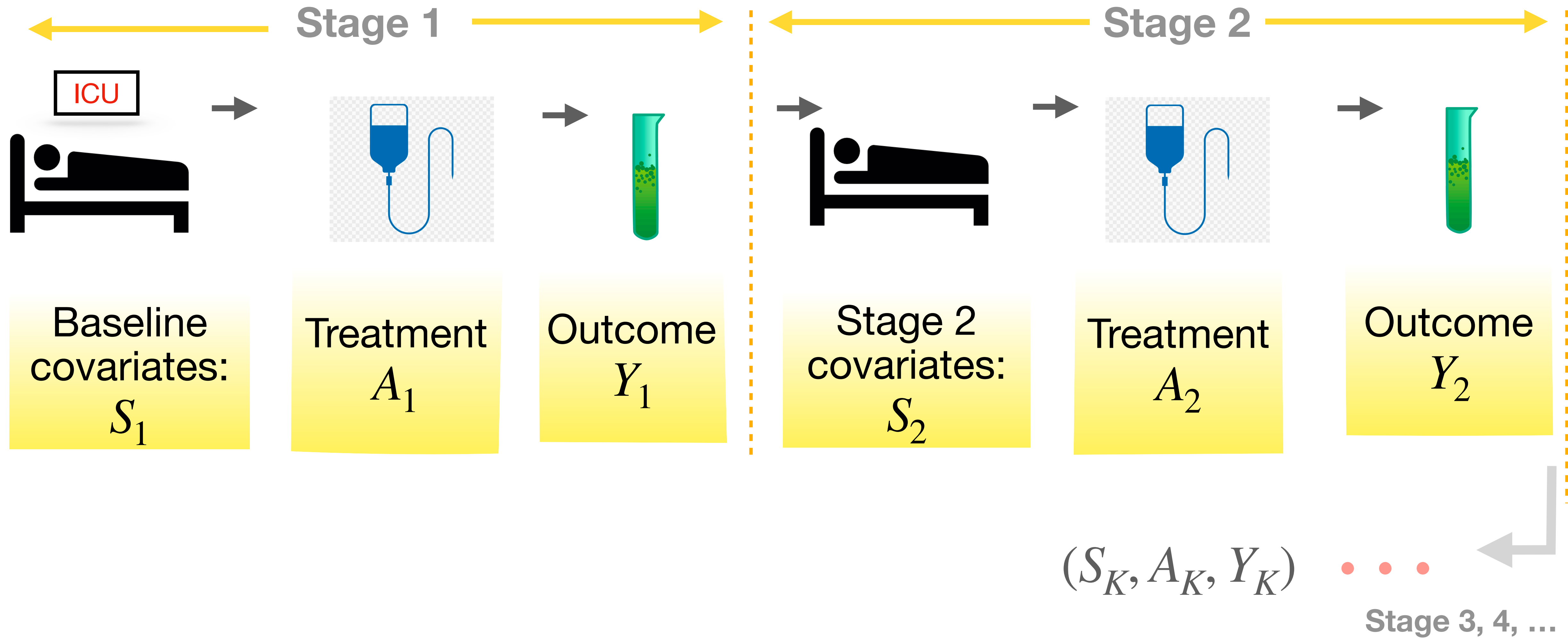




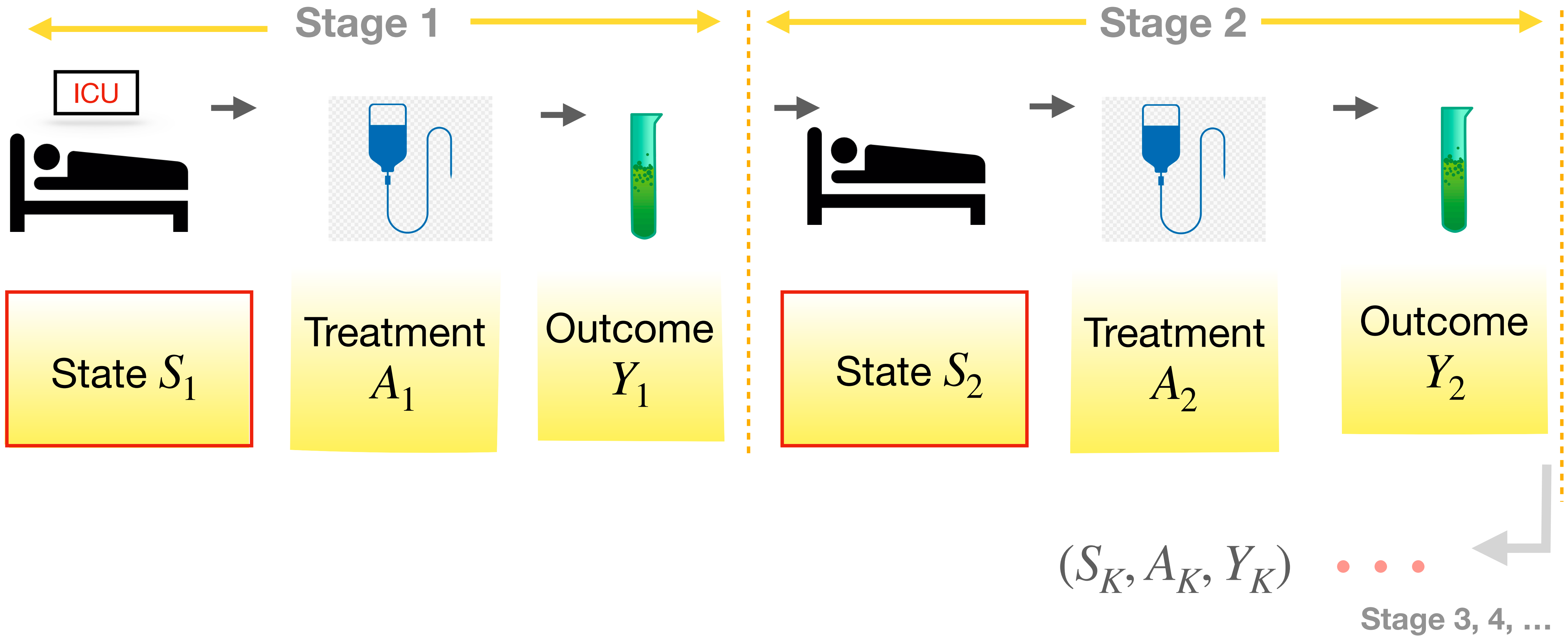
best treatment policy $\pi^* \implies \sum_{i=1}^K Y_i : \text{maximized}$

discharge
 (S_K, A_K, Y_K)
Stage 3, 4, ...

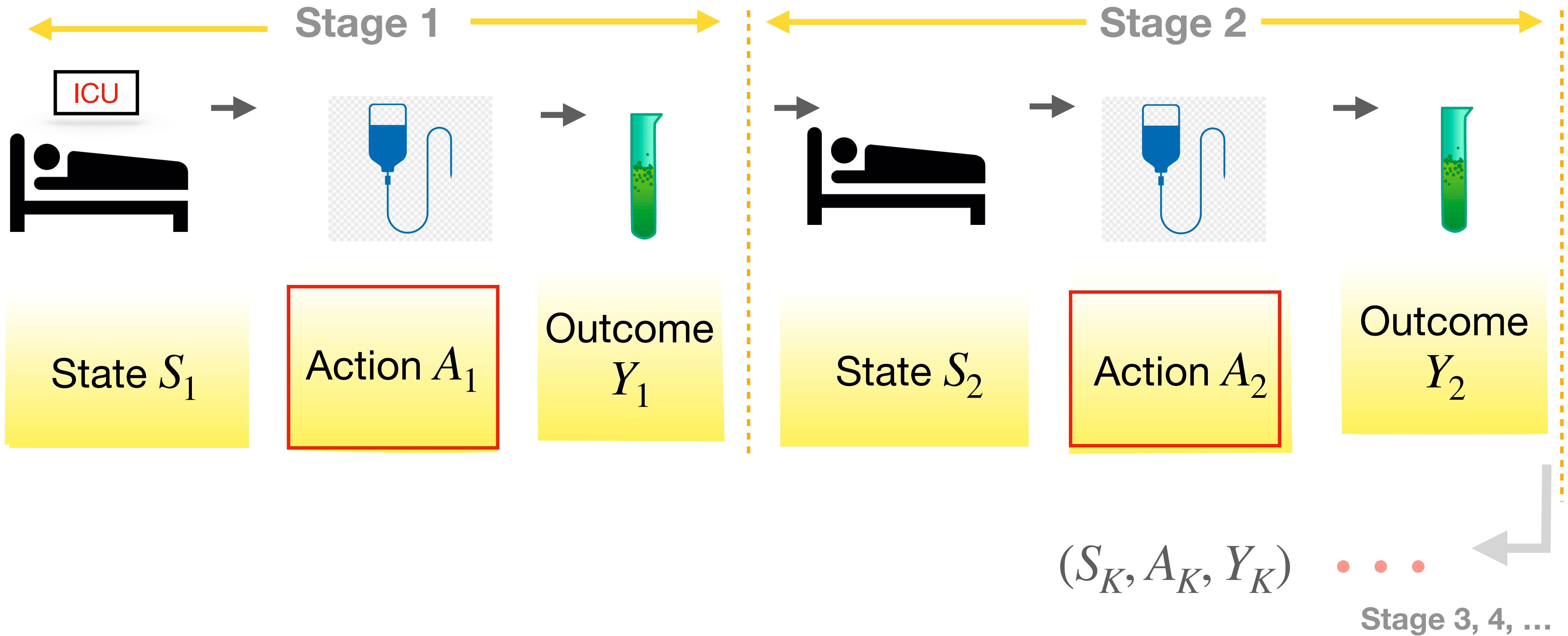
Offline Reinforcement Learning



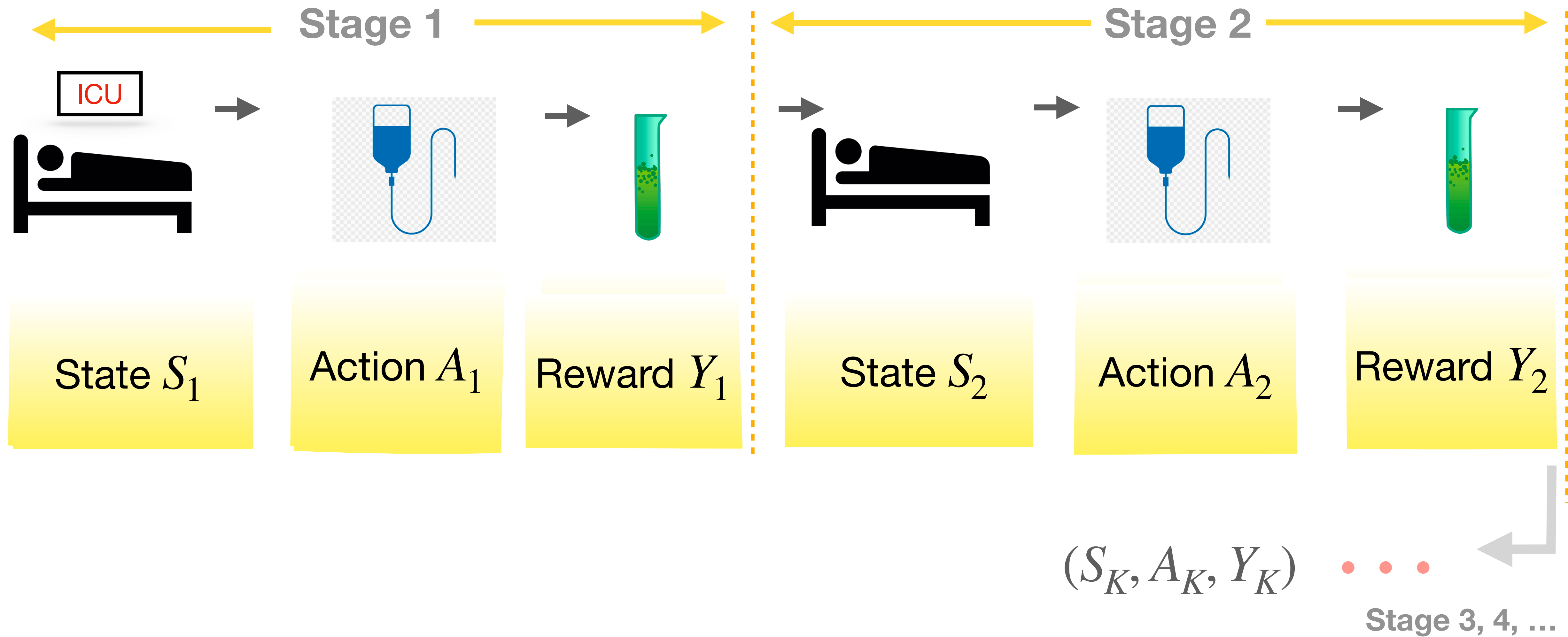
Offline Reinforcement Learning



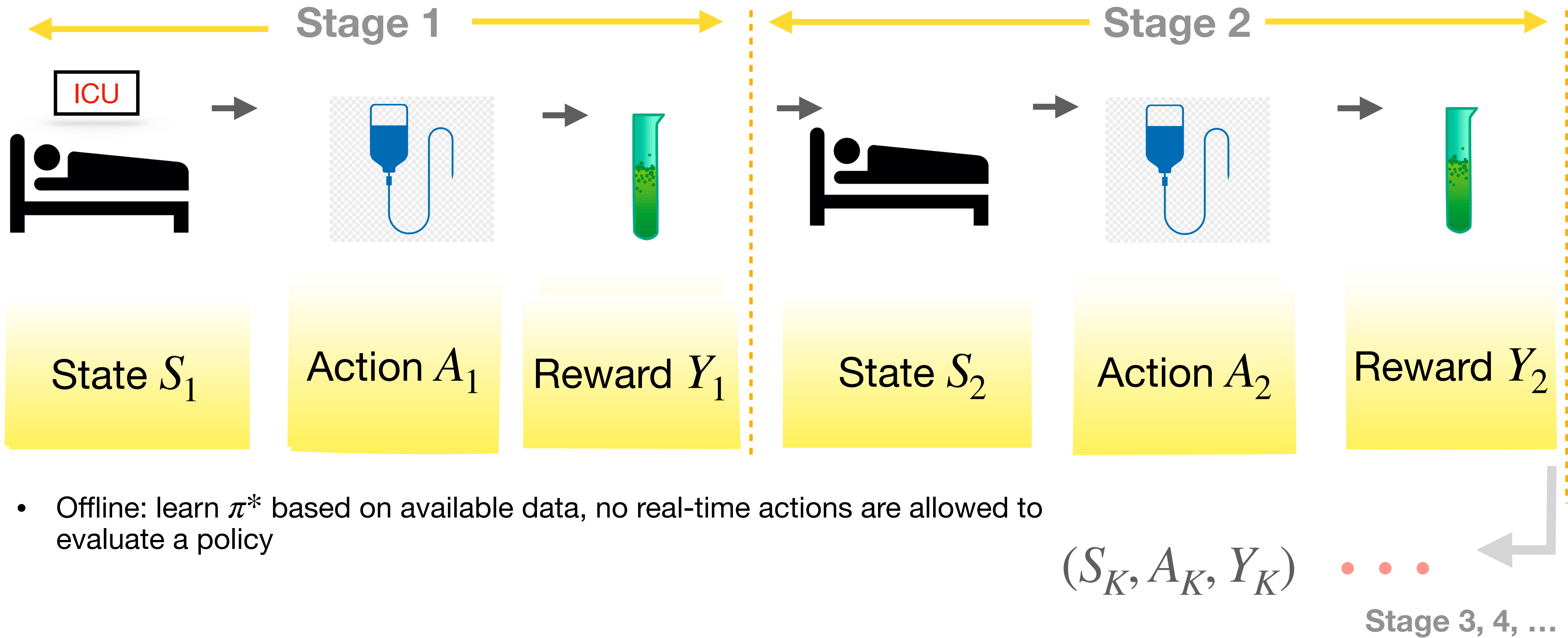
Offline Reinforcement Learning



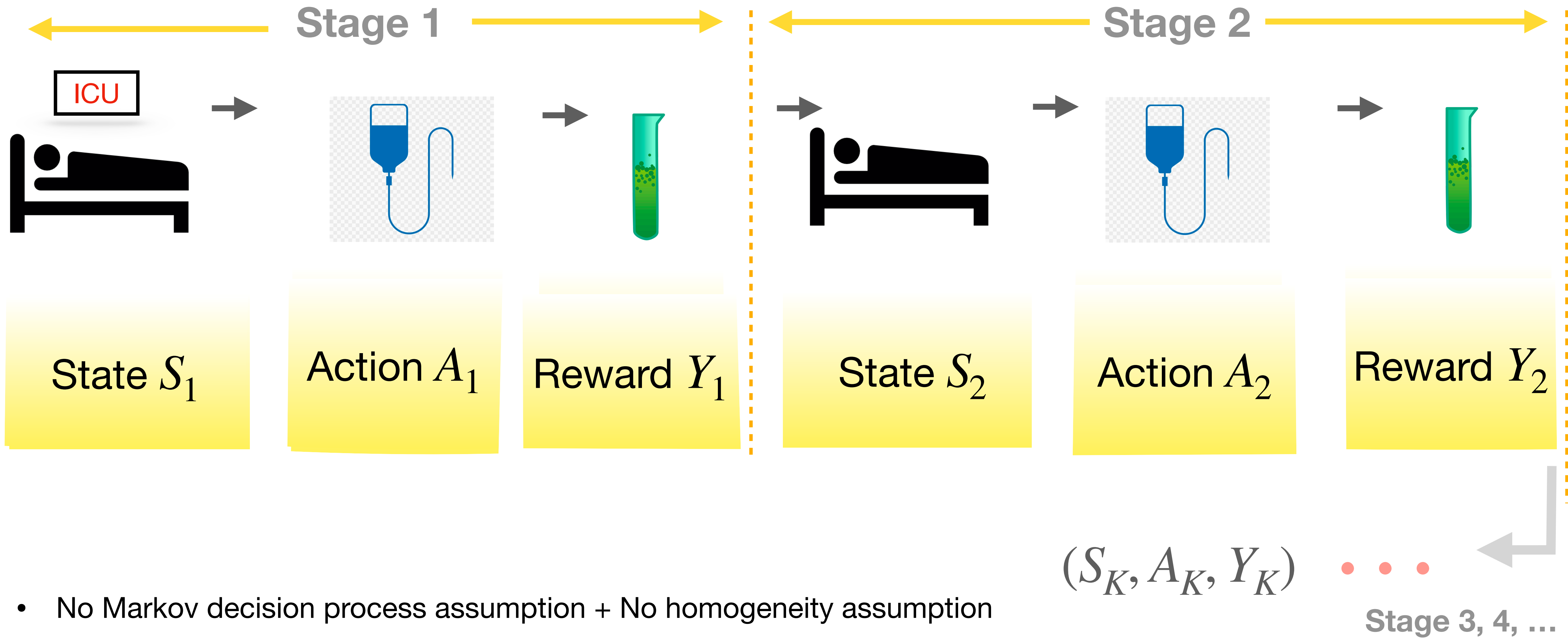
Offline Reinforcement Learning



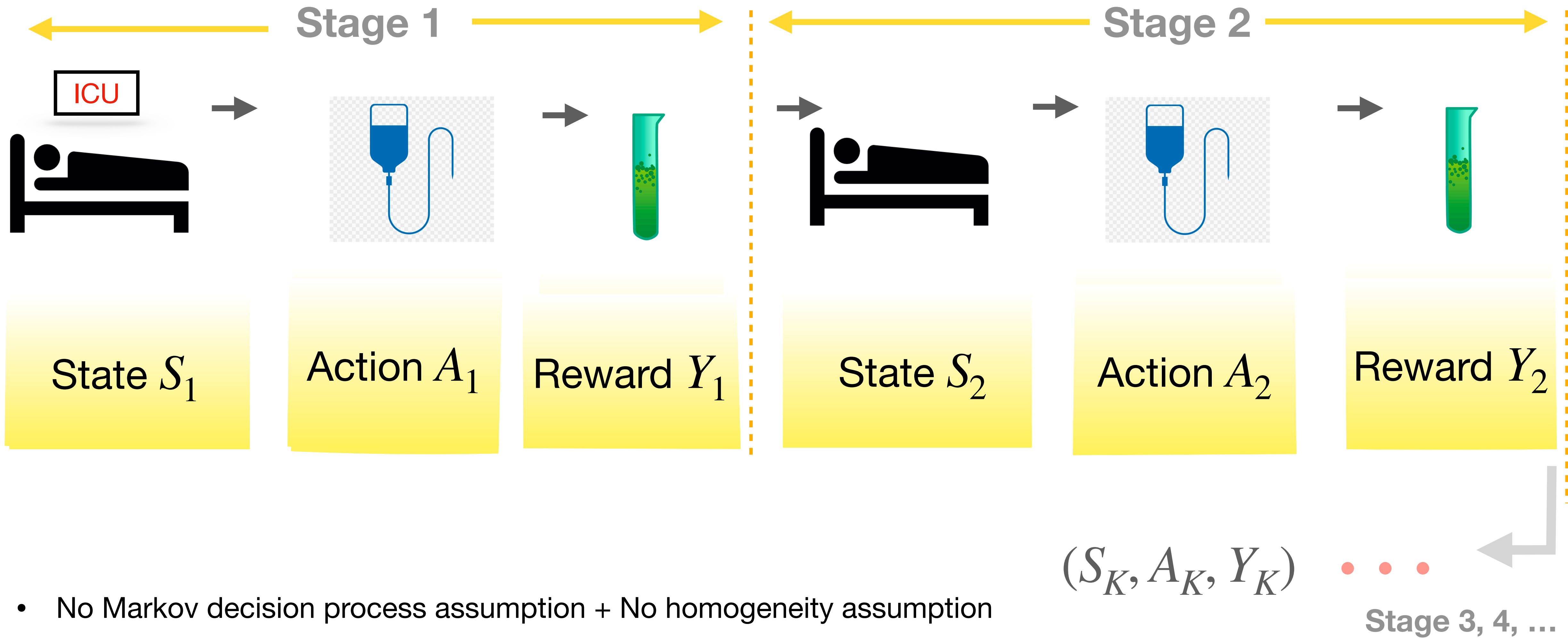
Offline Reinforcement Learning



Offline Reinforcement Learning



Offline Reinforcement Learning



- No Markov decision process assumption + No homogeneity assumption
- Hence called Full reinforcement learning

Outline

- Example: sepsis

- Problem formulation

A. Mathematical formulation

B. Existing approaches

- Proposed method
- Open questions

Outline

- Example: sepsis
- Problem formulation

A. Mathematical formulation

B. Existing approaches

- Proposed method
- Open questions

Mathematical formulation

History



History



H_1

First stage history

History



H_1

First stage history

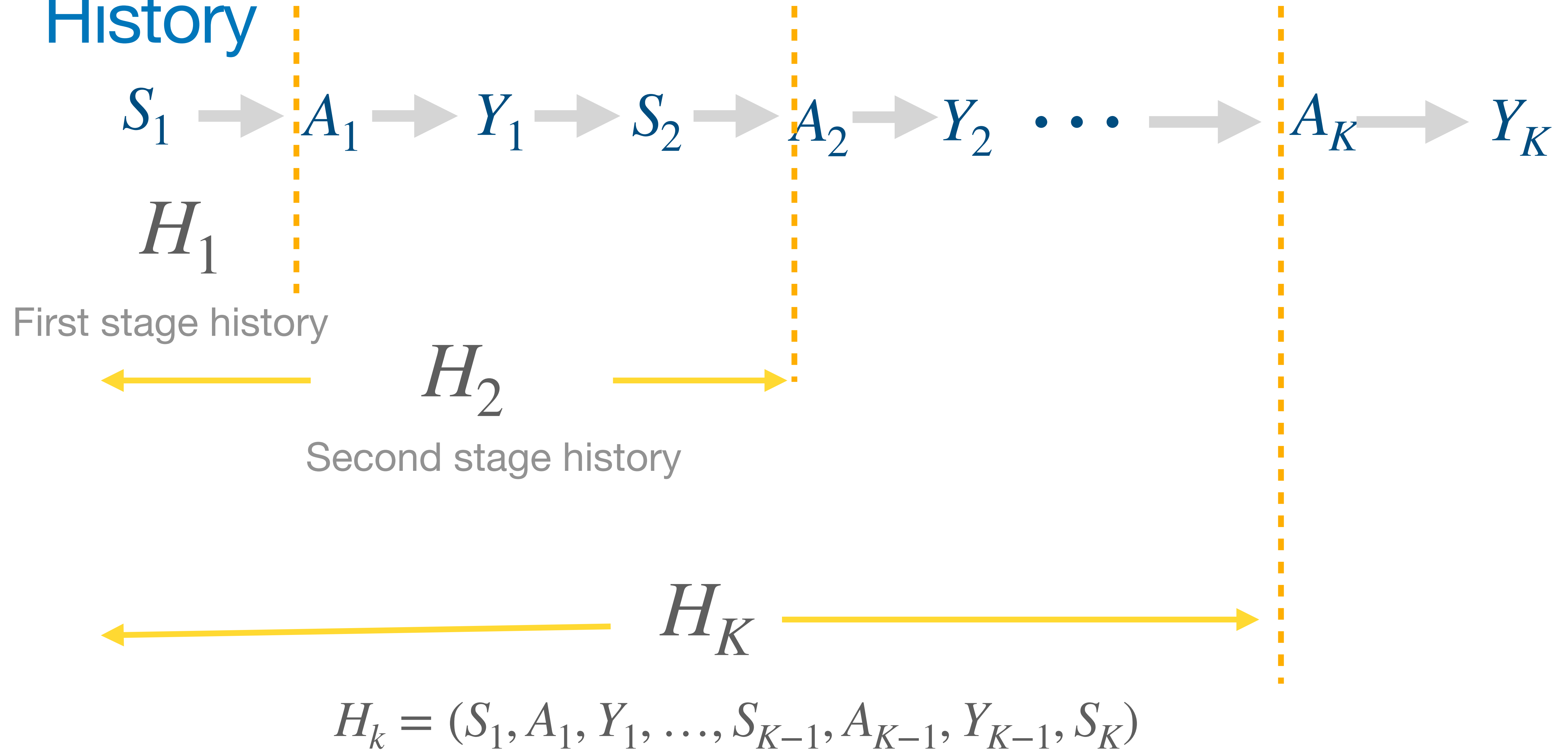


H_2

Second stage history



History



Policy

Policy

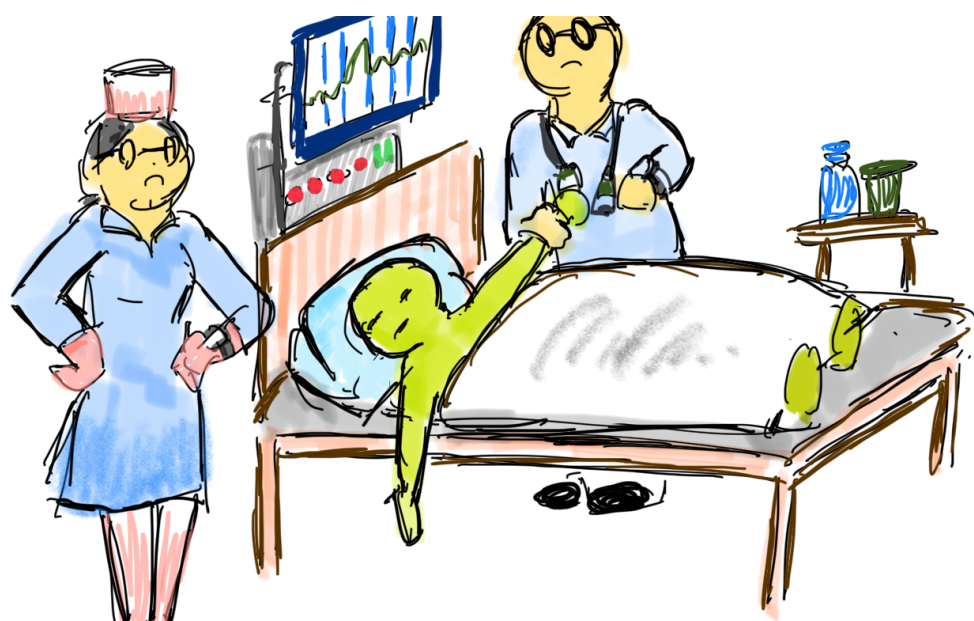
Stage 1



Look at H_1

Policy

Stage 1



Look at H_1

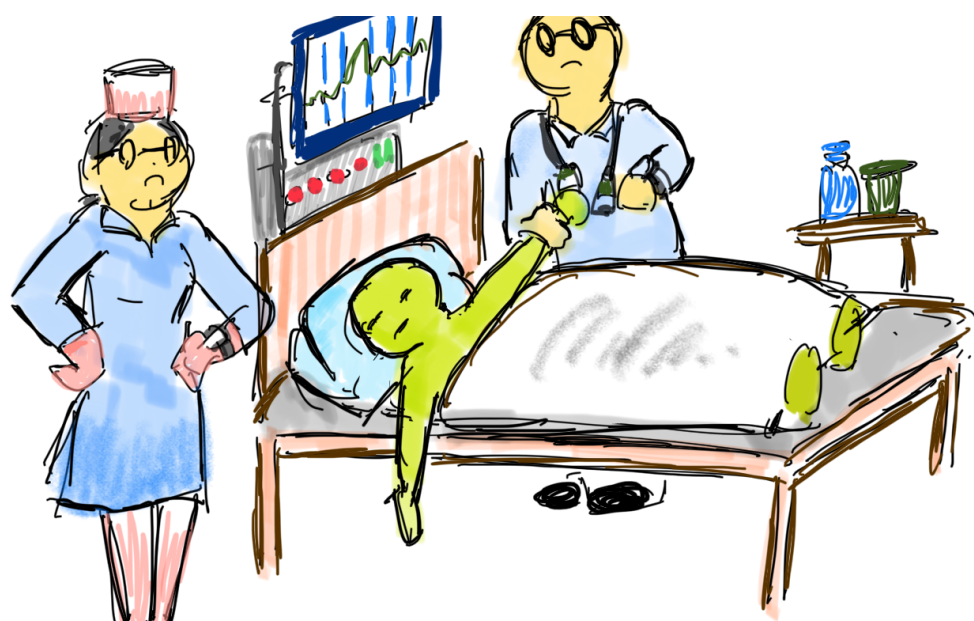


Choose IV-fluid
level

$$\pi_1(H_1) \in \mathcal{A} = \{\text{no fluid, low, mid, high}\}$$

Policy

Stage 1



Look at H_1



Choose IV-fluid
level

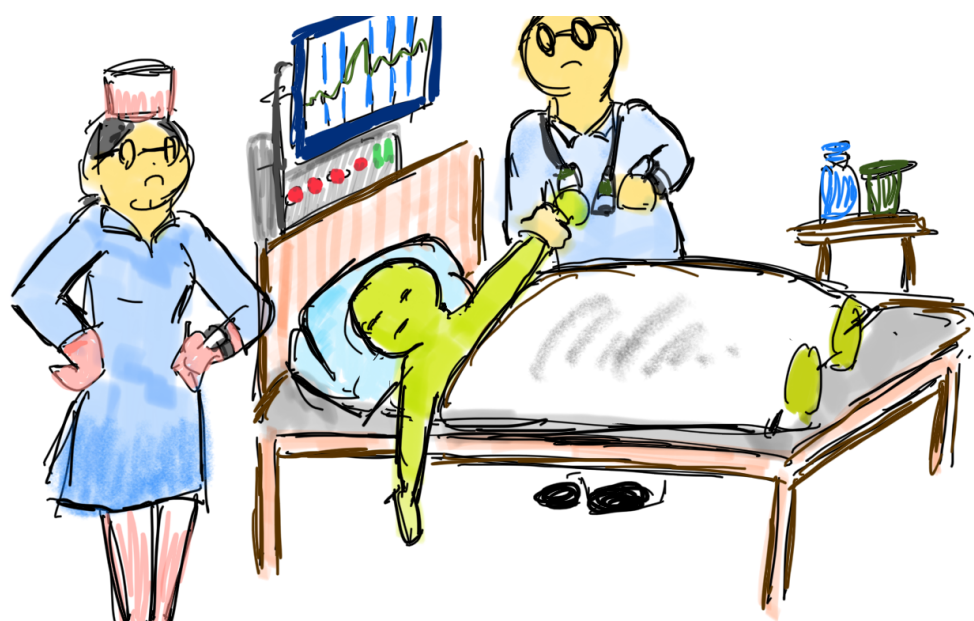
$$\pi_1(H_1) \in \mathcal{A} = \{\text{no fluid, low, mid, high}\}$$

Treatment Assignments

$$\pi_1 : H_1 \mapsto \mathcal{A}$$

Policy

Stage 1



Look at H_1



Choose IV-fluid
level

$$\pi_1(H_1) \in \mathcal{A} = \{\text{no fluid, low, mid, high}\}$$

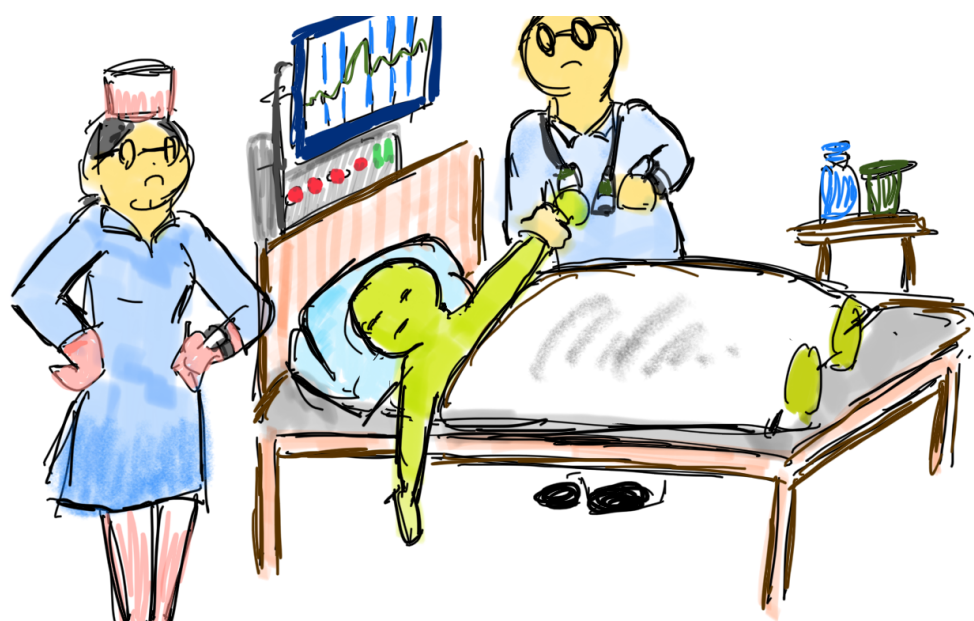
Stages 2, 3, 4, ...

Treatment Assignments

$$\pi_1 : H_1 \mapsto \mathcal{A}$$

Policy

Stage 1



Look at H_1



Choose IV-fluid
level

$$\pi_1(H_1) \in \mathcal{A} = \{\text{no fluid, low, mid, high}\}$$

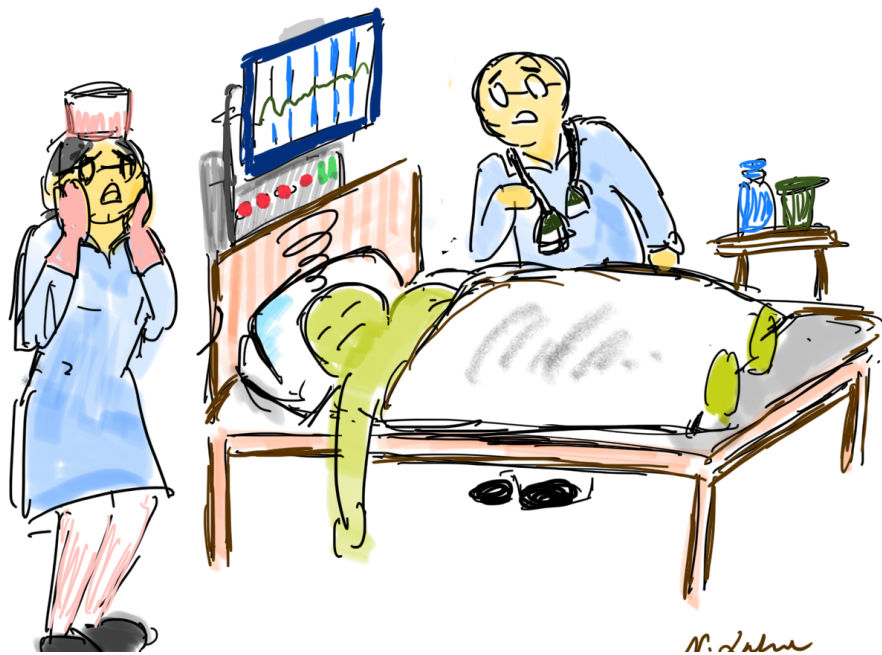
Treatment Assignments

$$\pi_1 : H_1 \mapsto \mathcal{A}$$

Stages 2, 3, 4, ...



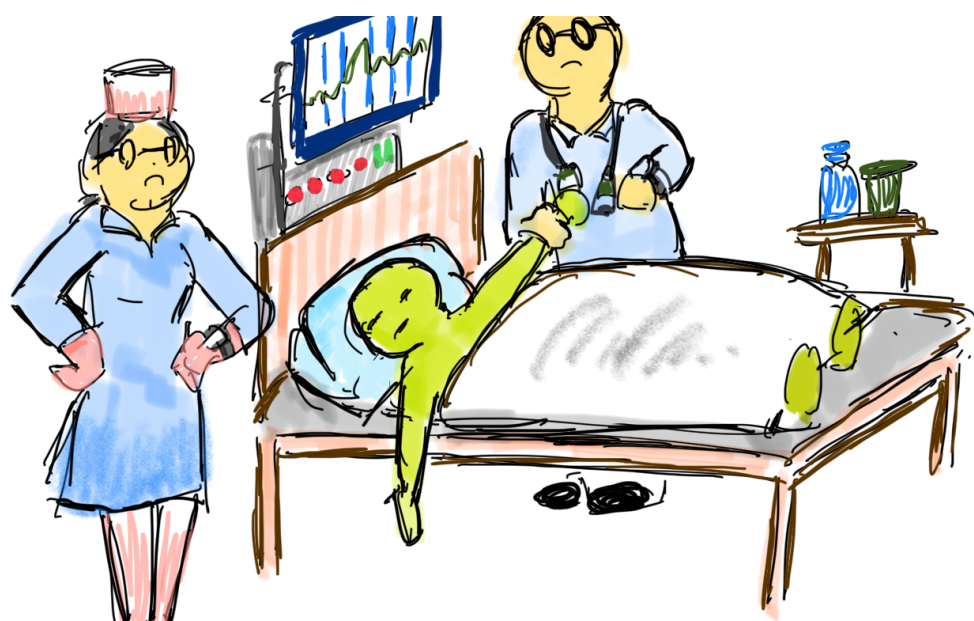
Stage k



Look at H_k

Policy

Stage 1



Look at H_1



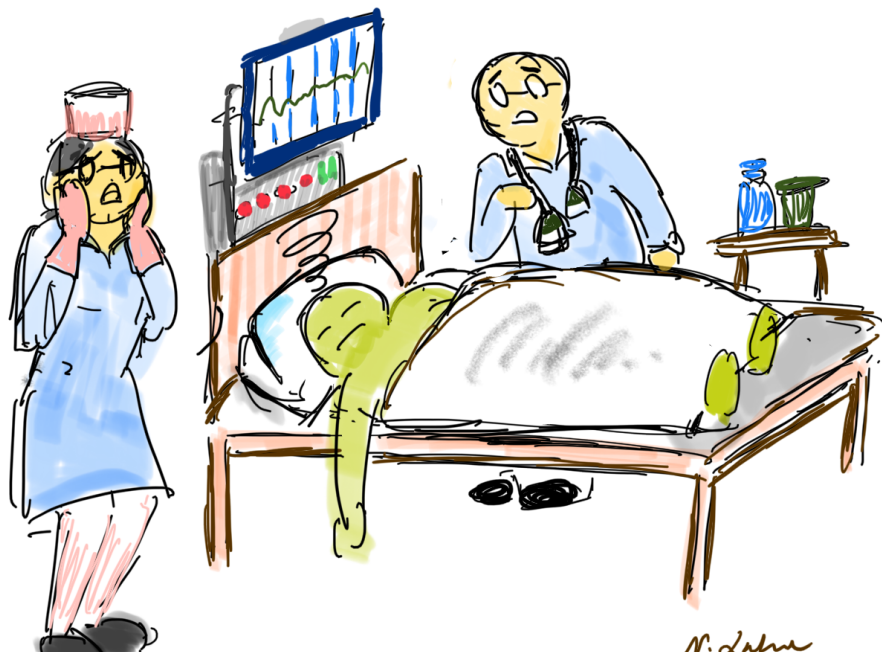
Choose IV-fluid
level

$$\pi_1(H_1) \in \mathcal{A} = \{\text{no fluid, low, mid, high}\}$$

Stages 2, 3, 4, ...



Stage k



Look at H_k



Choose IV-fluid
level

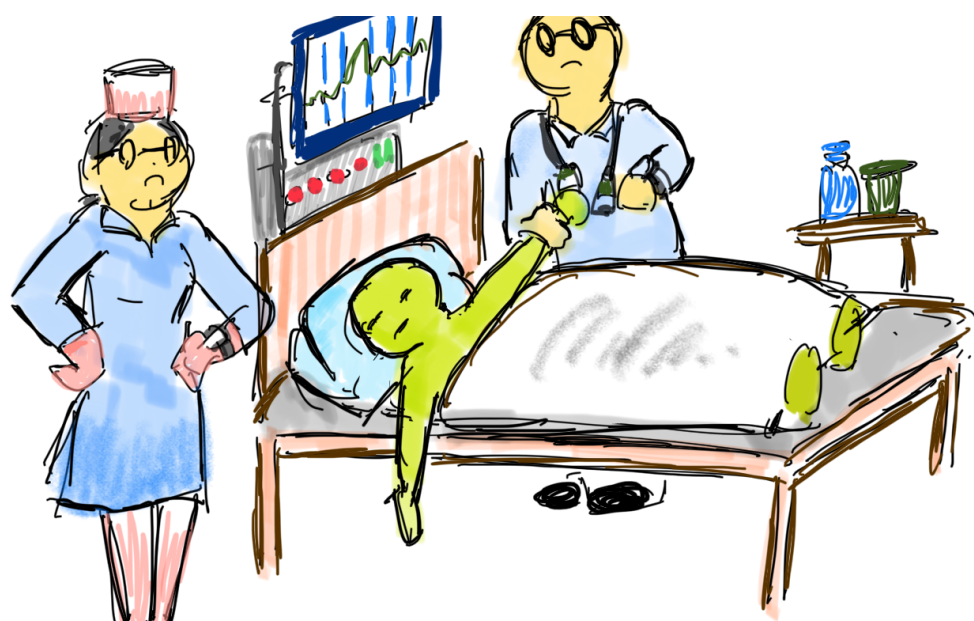
$$\pi_k(H_k) \in \{\text{no fluid, low, mid, high}\}$$

Treatment Assignments

$$\pi_1 : H_1 \mapsto \mathcal{A}$$

Policy

Stage 1



Look at H_1



Choose IV-fluid
level

$$\pi_1(H_1) \in \mathcal{A} = \{\text{no fluid, low, mid, high}\}$$

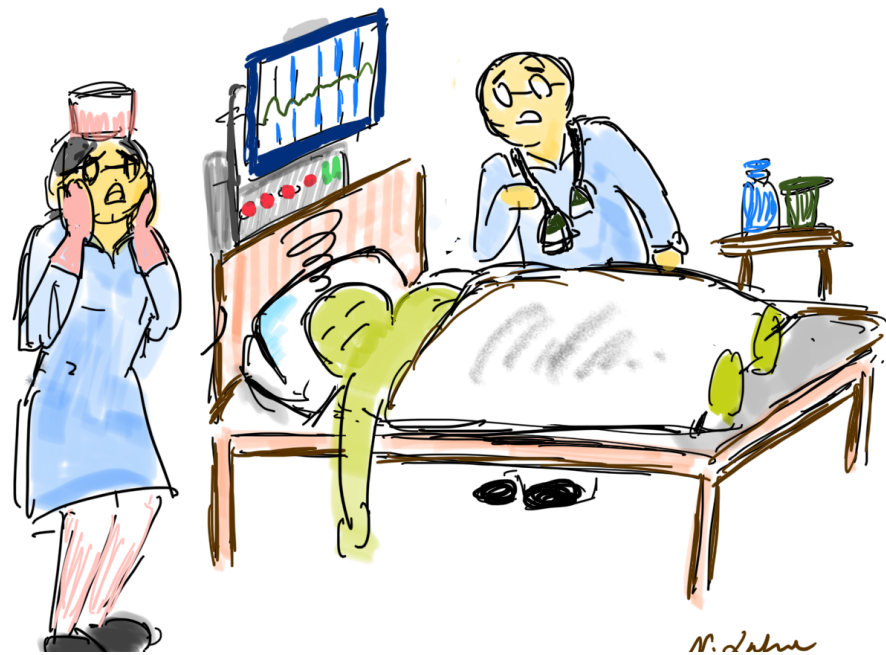
Treatment Assignments

$$\pi_1 : H_1 \mapsto \mathcal{A}$$

Stages 2, 3, 4, ...



Stage k



Look at H_k



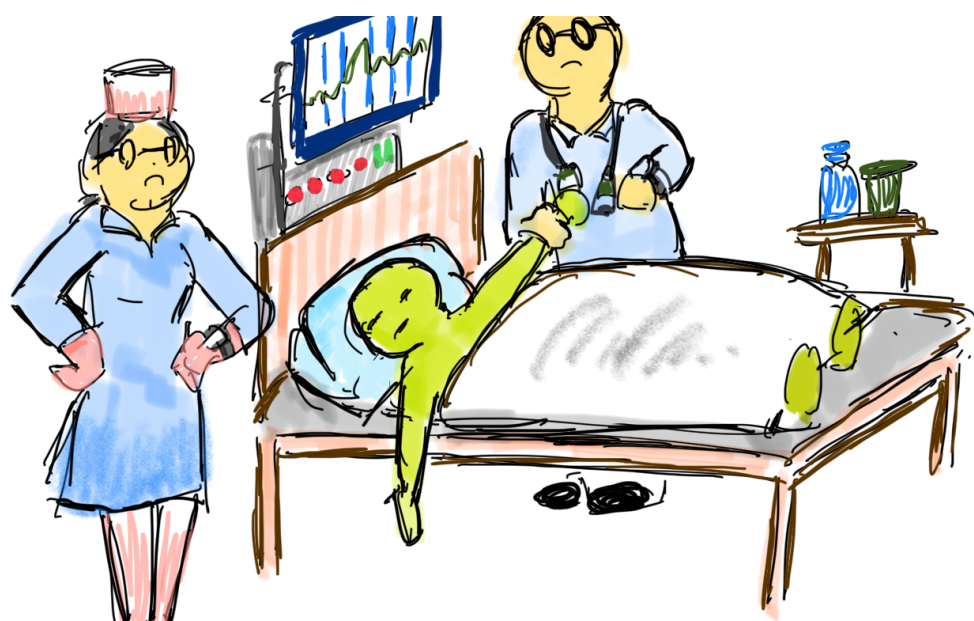
Choose IV-fluid
level

$$\pi_k(H_k) \in \{\text{no fluid, low, mid, high}\}$$

$$\pi_k : H_k \mapsto \mathcal{A}$$

Policy

Stage 1



Look at H_1



Choose IV-fluid
level

$$\pi_1(H_1) \in \mathcal{A} = \{\text{no fluid, low, mid, high}\}$$

Treatment Assignments

$$\pi_1 : H_1 \mapsto \mathcal{A}$$

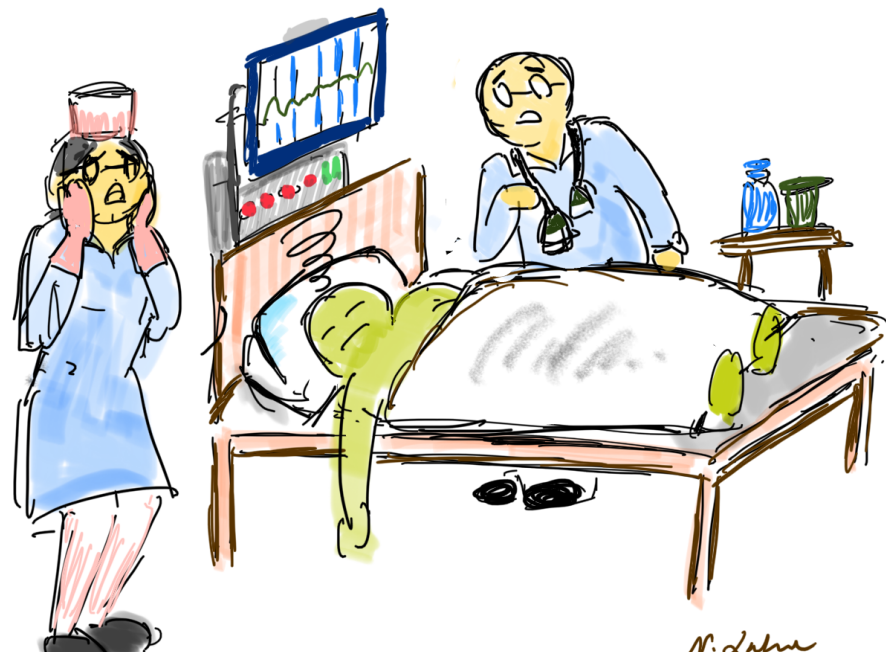
Stages 2, 3, 4, ...



Policy

$$\pi = (\pi_1, \dots, \pi_K)$$

Stage k



Look at H_k



Choose IV-fluid
level

$$\pi_k(H_k) \in \{\text{no fluid, low, mid, high}\}$$

$$\pi_k : H_k \mapsto \mathcal{A}$$

Example of policy in full RL

Example of policy in full RL

Stage 1:

Example of policy in full RL

Stage 1:

$$\pi_1$$

Example of policy in full RL

Systolic blood pressure ≤ 90 mm
Hg

Stage 1:

π_1

Example of policy in full RL

Stage 1:

π_1

Systolic blood pressure ≤ 90 mm
Hg

Yes



Example of policy in full RL

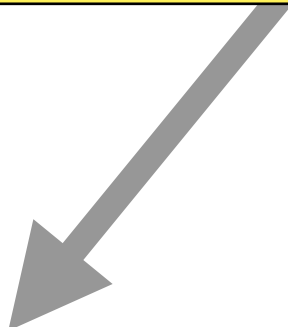
Stage 1:

π_1

IV

Systolic blood pressure \leq 90 mm
Hg

Yes

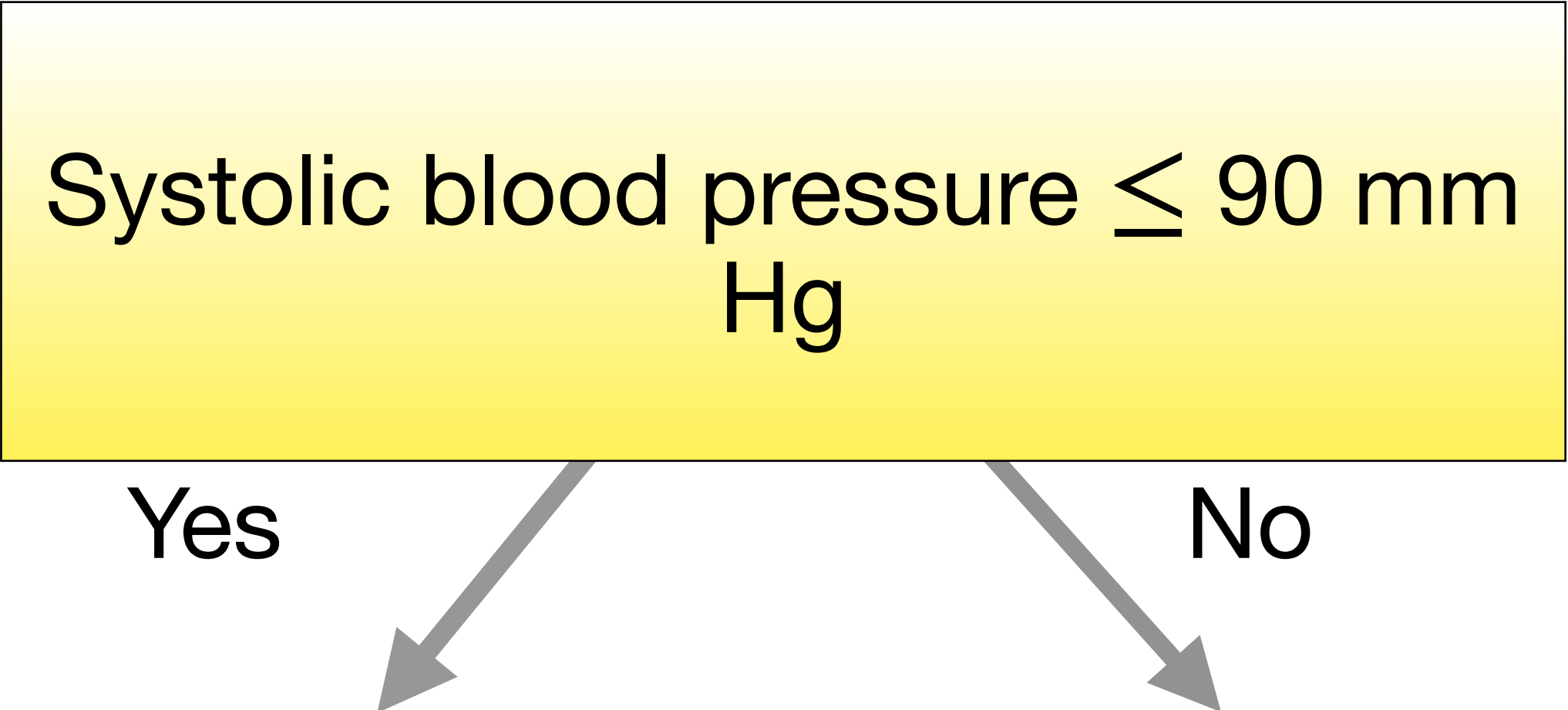


Example of policy in full RL

Stage 1:

π_1

IV



Example of policy in full RL

Stage 1:

π_1

Systolic blood pressure \leq 90 mm
Hg

Yes

No

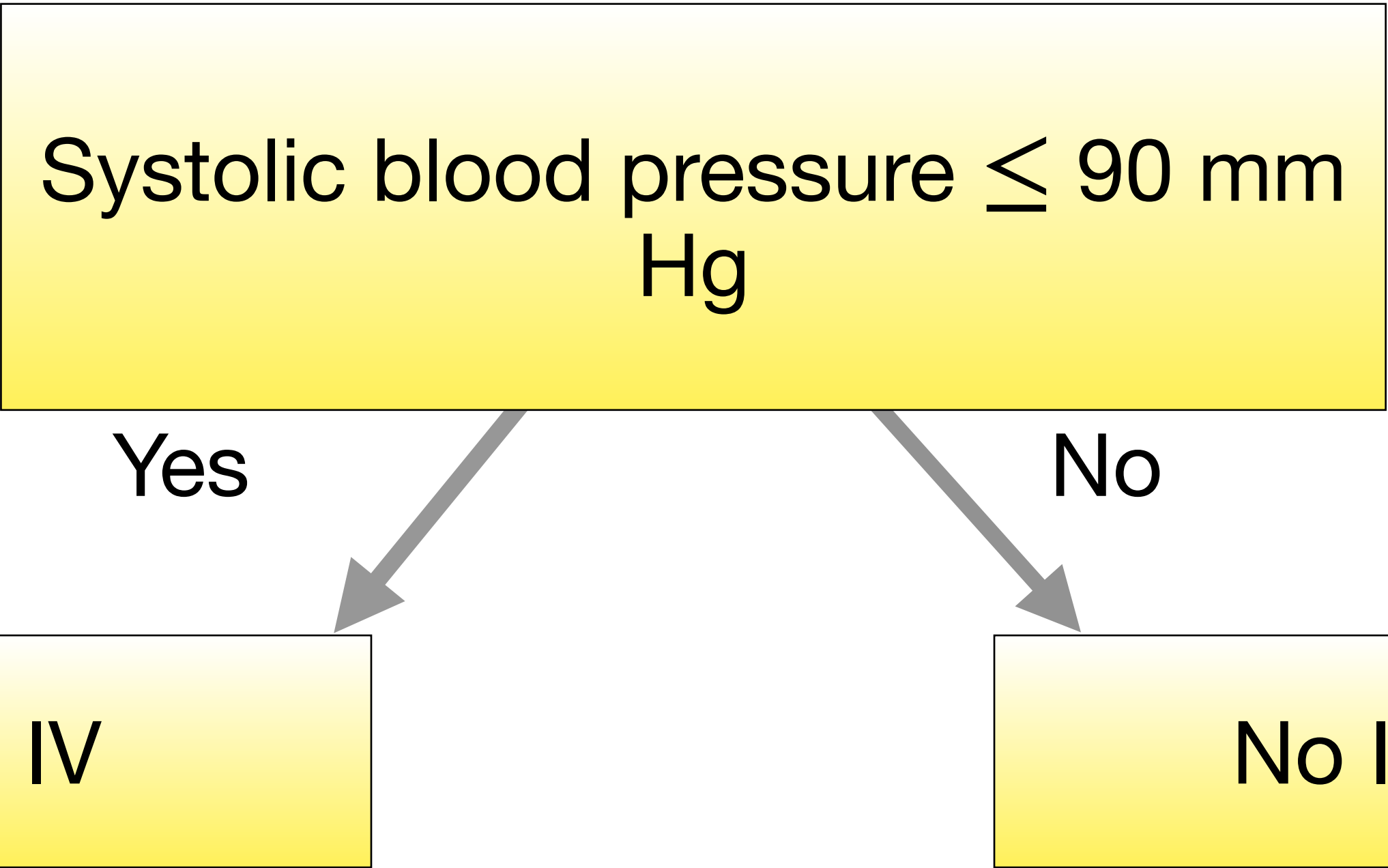
IV

No IV

Example of policy in full RL

Stage 1:

π_1



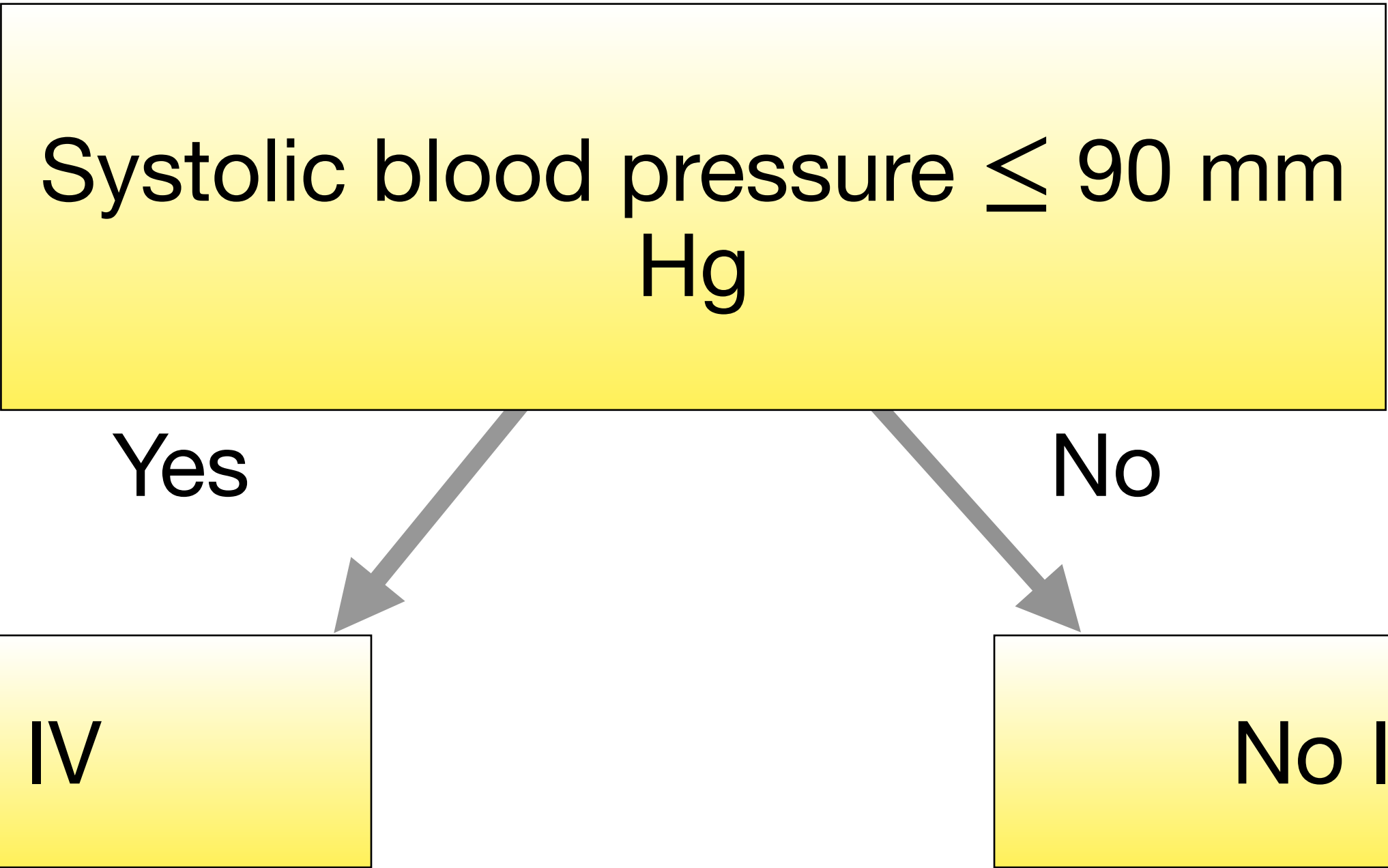
Stage k :

$t = 2, \dots, K.$

Example of policy in full RL

Stage 1:

π_1



Stage k :

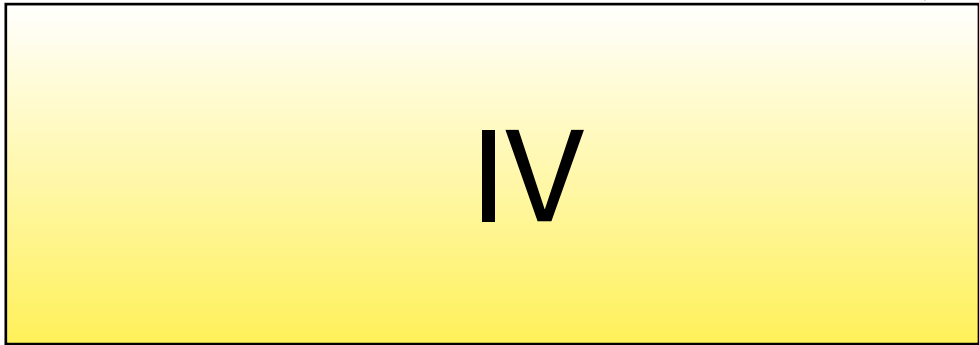
$t = 2, \dots, K.$

π_k

Example of policy in full RL

Stage 1:

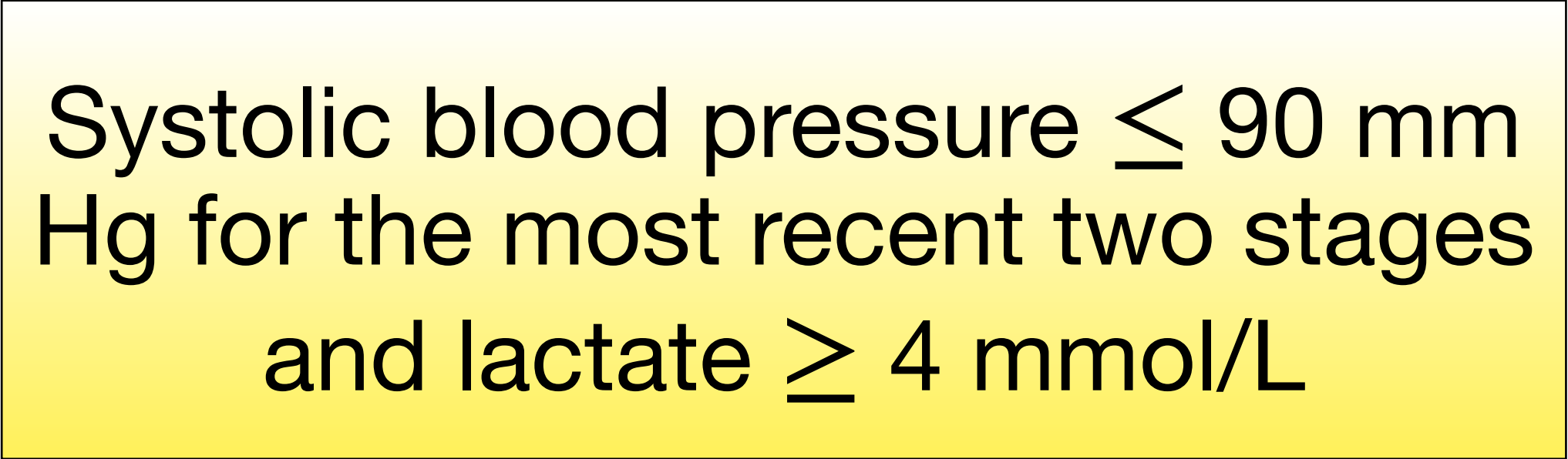
π_1



Stage k :

$t = 2, \dots, K.$

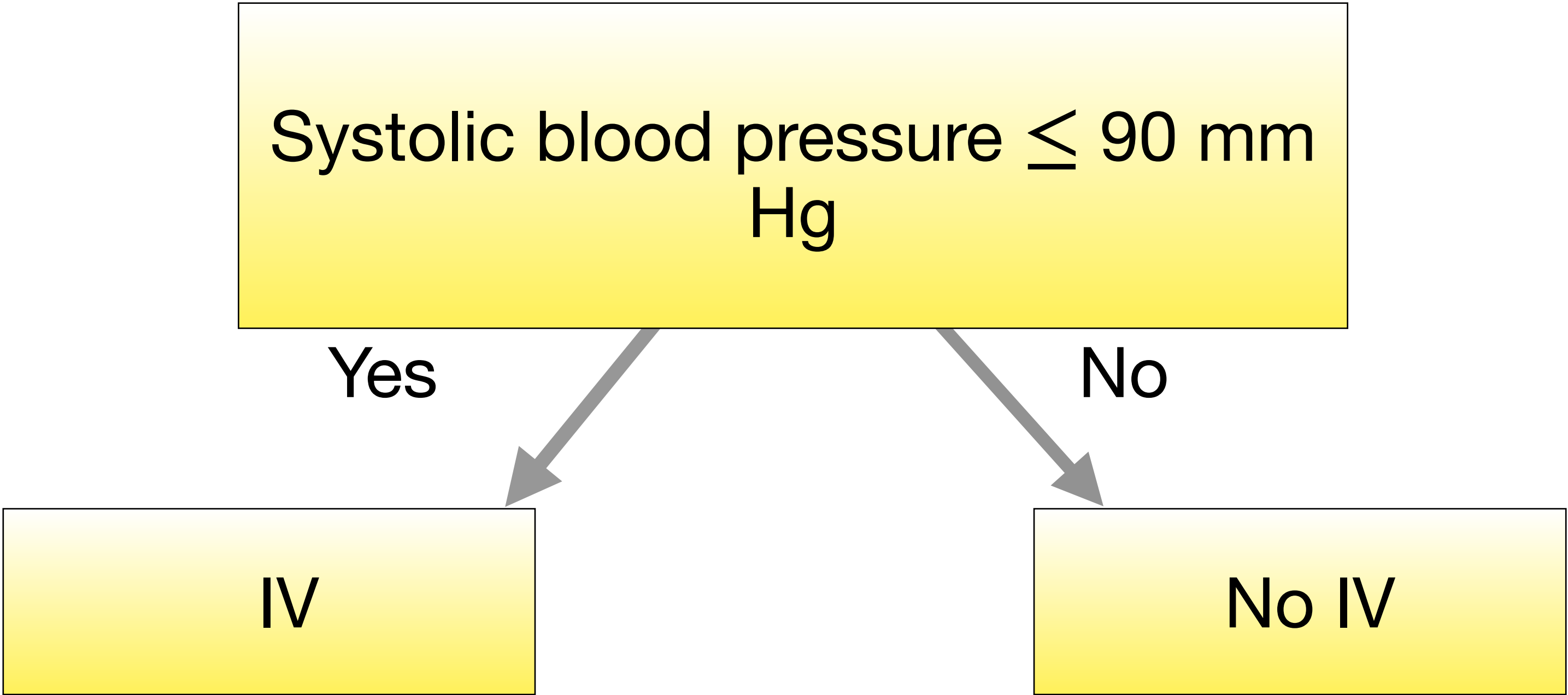
π_k



Example of policy in full RL

Stage 1:

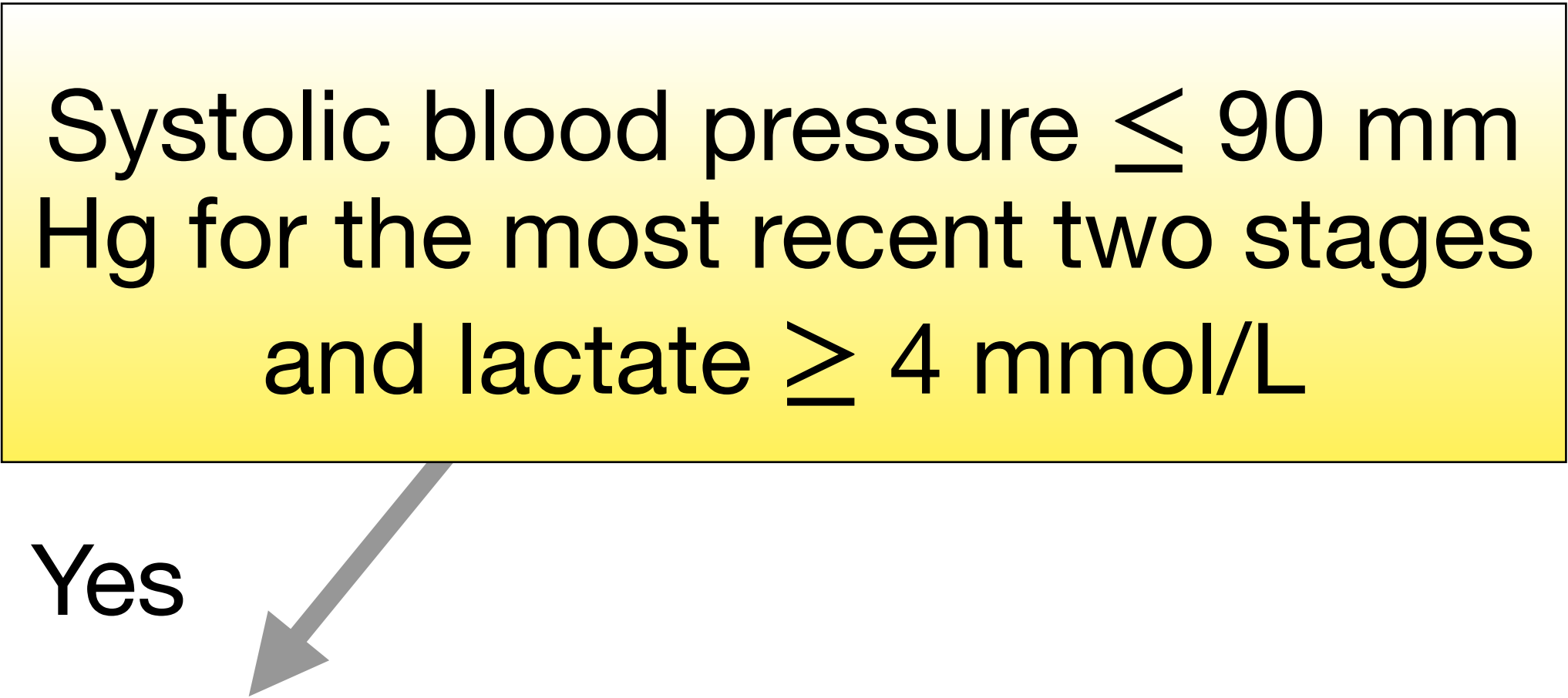
π_1



Stage k :

$t = 2, \dots, K.$

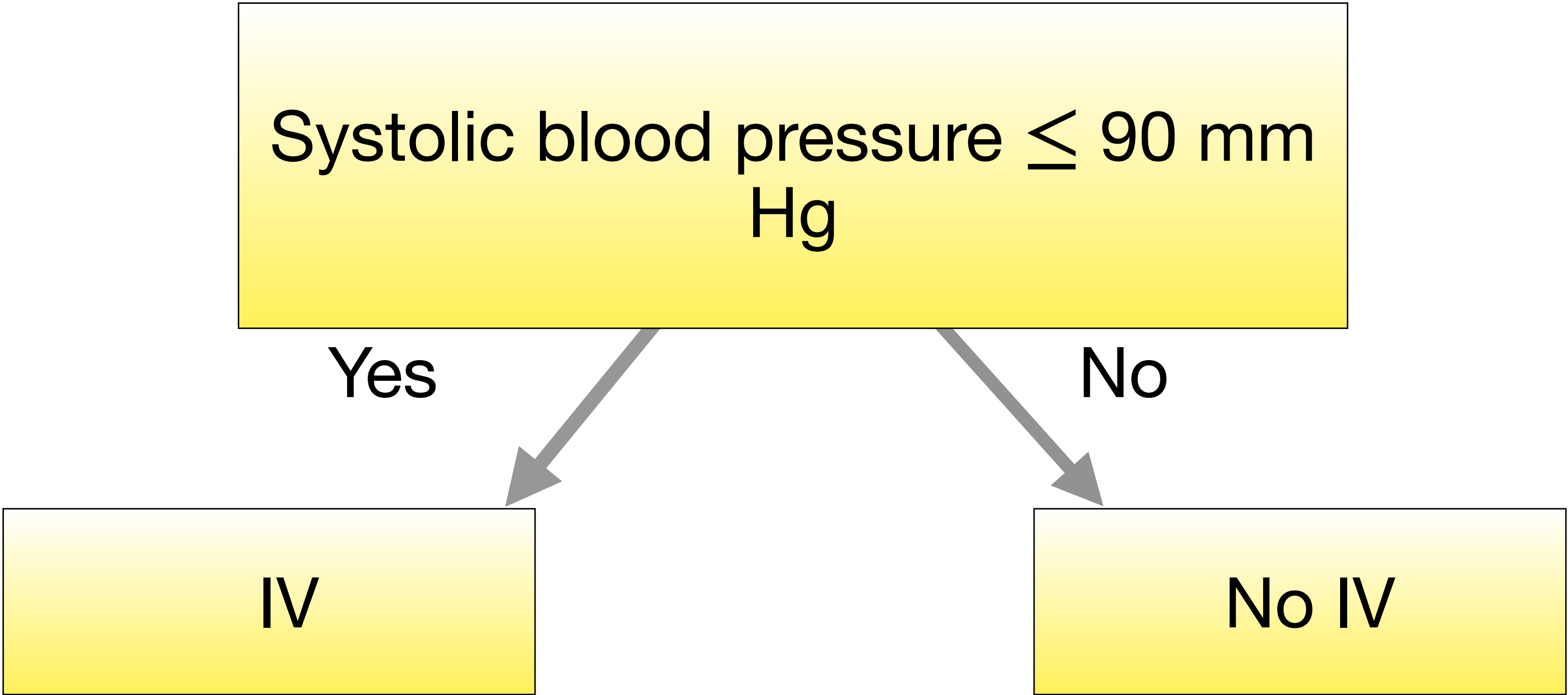
π_k



Example of policy in full RL

Stage 1:

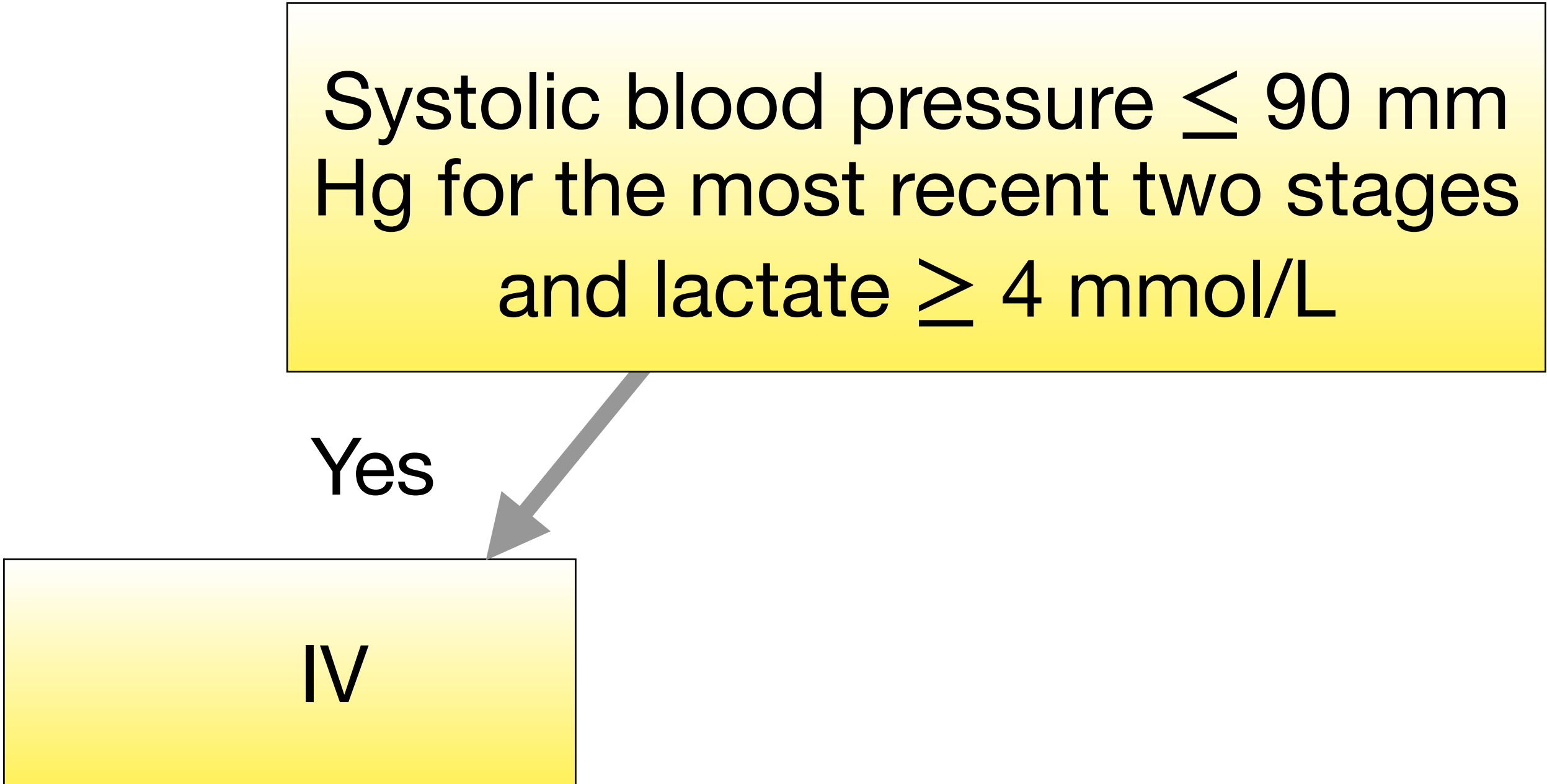
π_1



Stage k :

$t = 2, \dots, K.$

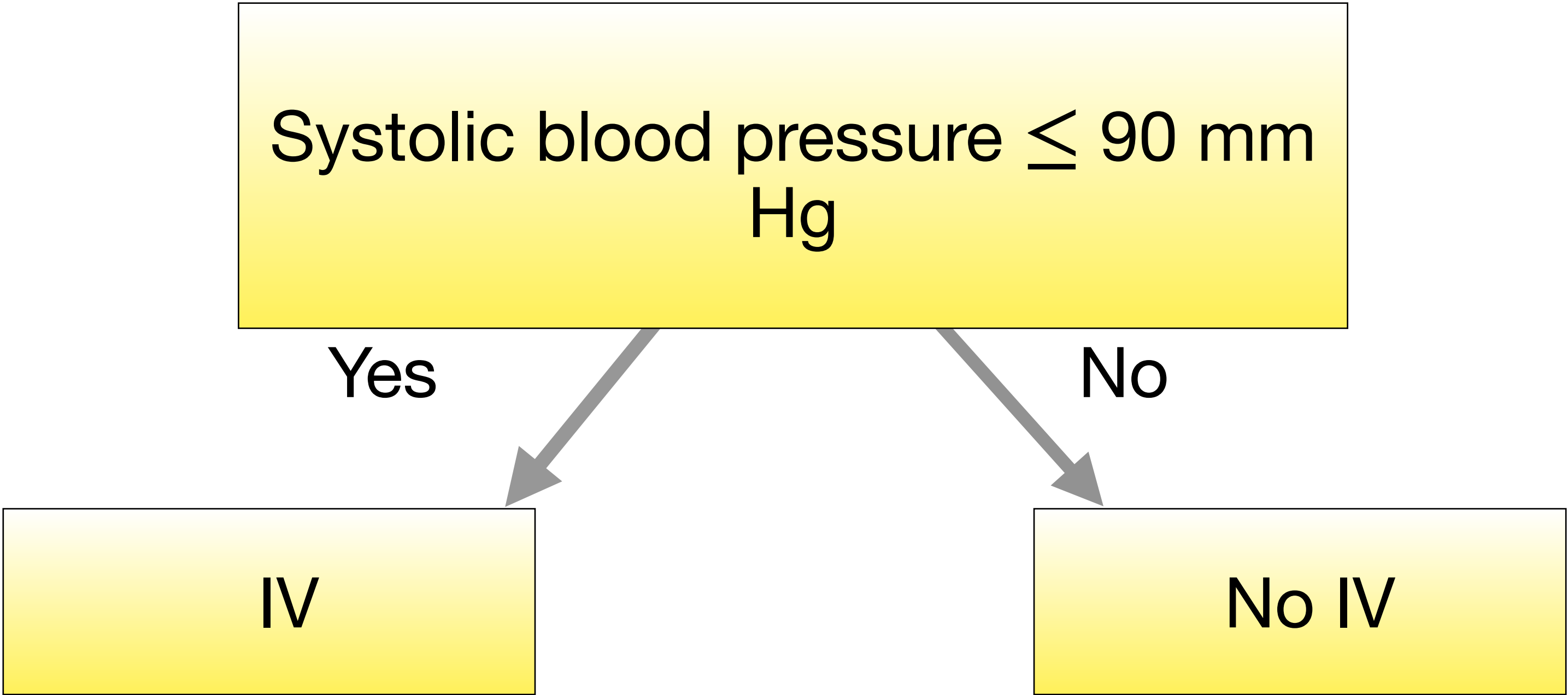
π_k



Example of policy in full RL

Stage 1:

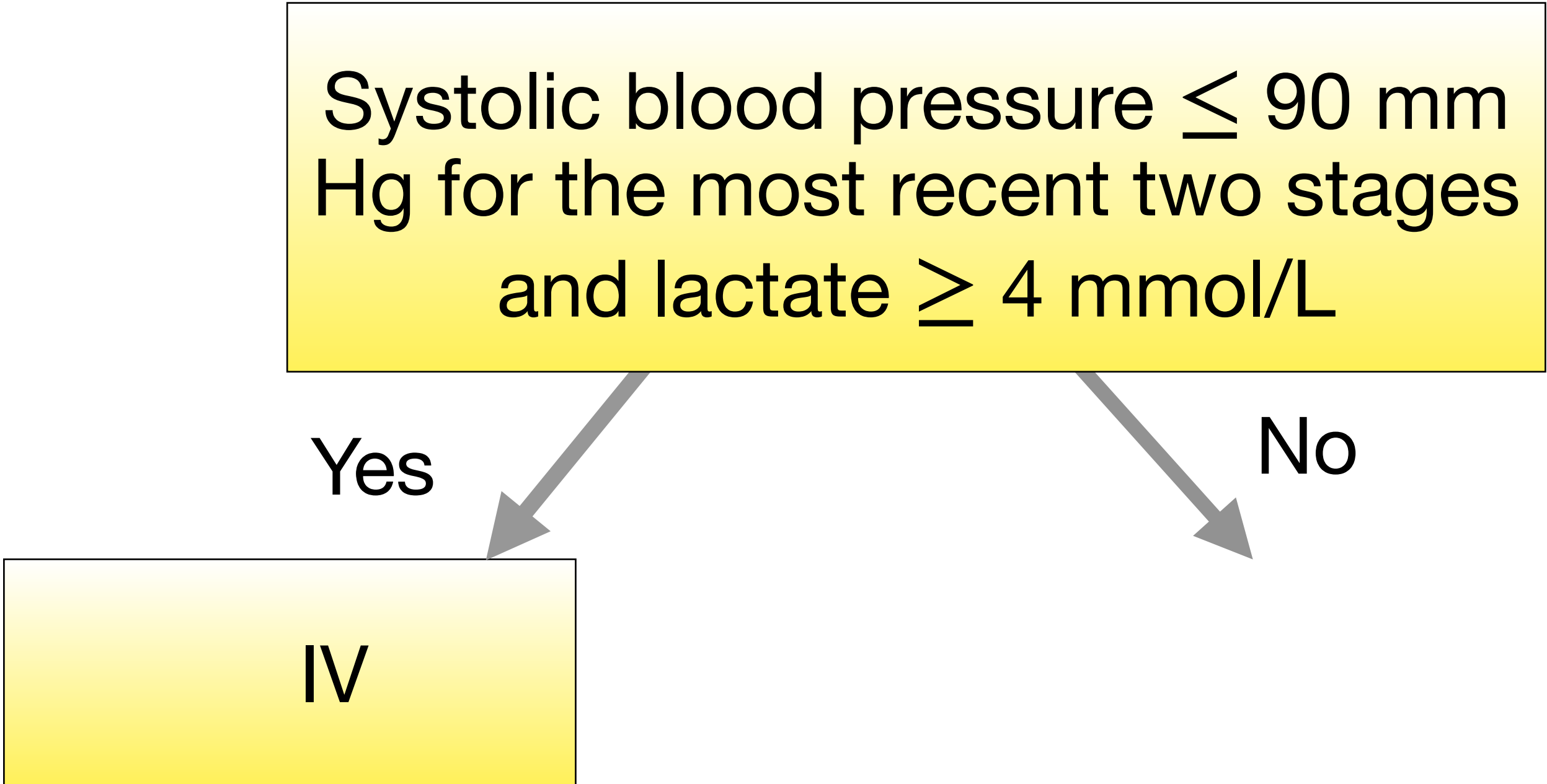
π_1



Stage k :

$t = 2, \dots, K.$

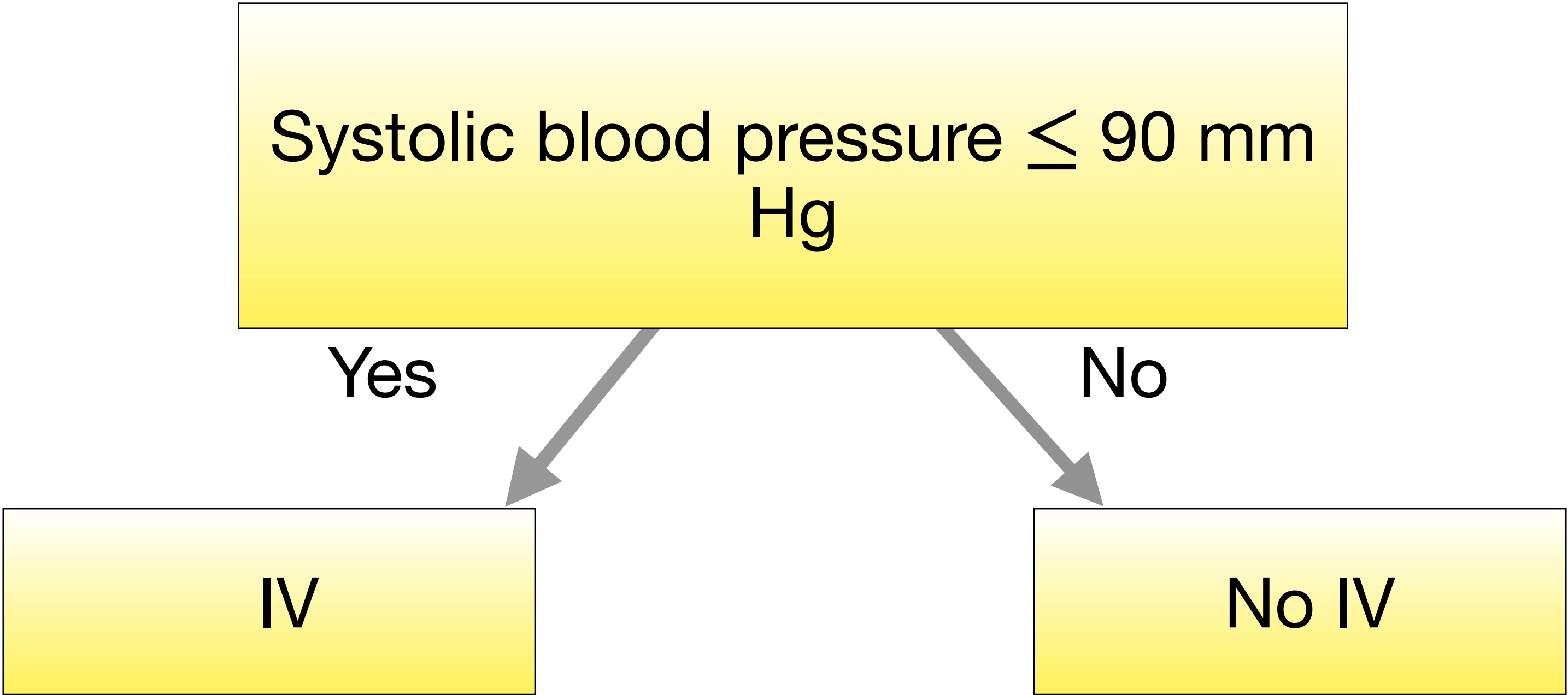
π_k



Example of policy in full RL

Stage 1:

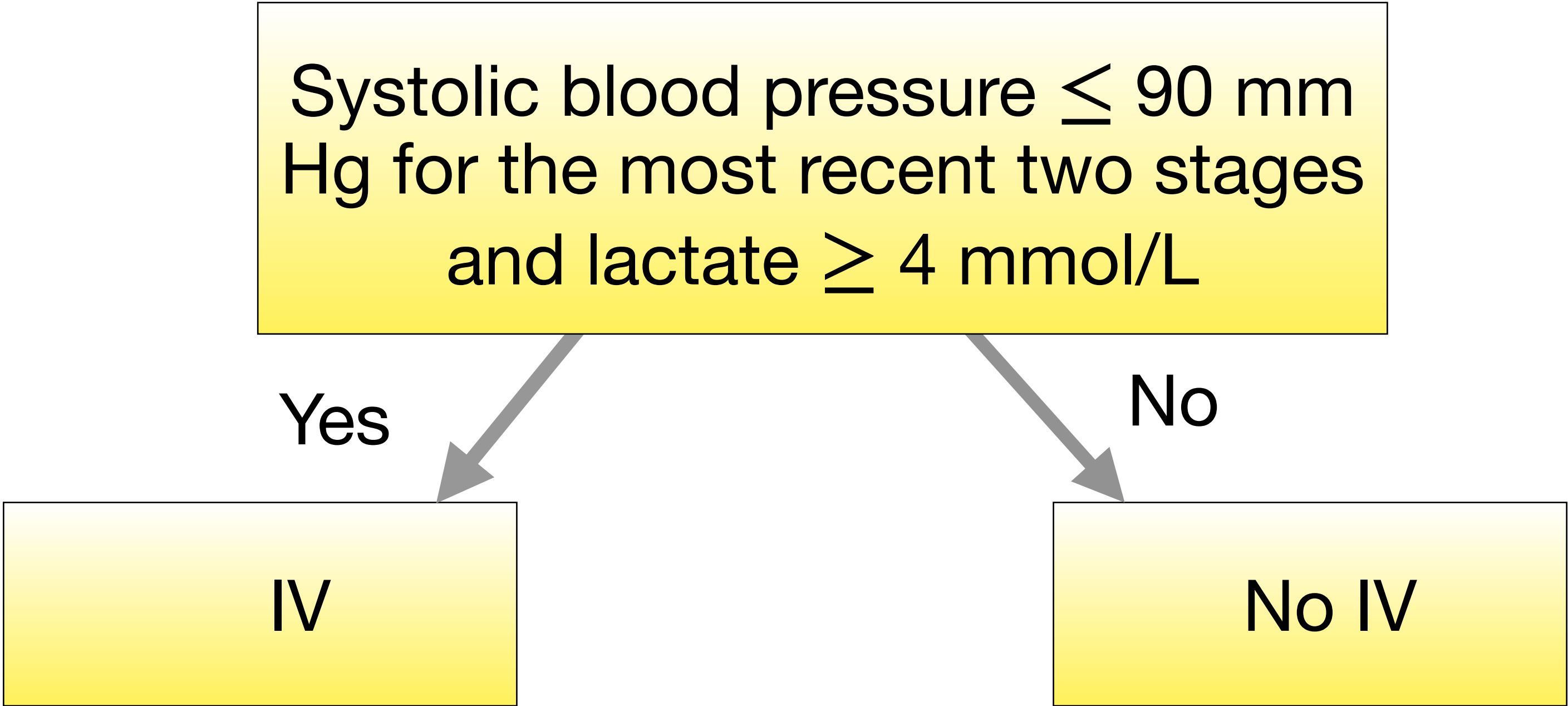
π_1



Stage k :

$t = 2, \dots, K.$

π_k



Optimal treatment policy: value function

Optimal treatment policy: value function

$Y_k(\pi)$: potential outcome at stage k had policy π been followed

Optimal treatment policy: value function

$Y_k(\pi)$: potential outcome at stage k had policy π been followed

Value function of π :

$$V^\pi = E[Y_1(\pi) \dots + Y_K(\pi)]$$

Optimal treatment policy: value function

$Y_k(\pi)$: potential outcome at stage k had policy π been followed

Value function of π :

$$V^\pi = E[Y_1(\pi) \dots + Y_K(\pi)]$$

$$\pi^* = \operatorname{argmax}_{\pi} V^\pi$$

Optimal treatment policy: value function

$Y_k(\pi)$: potential outcome at stage k had policy π been followed

Value function of π :

$$V^\pi = E[Y_1(\pi) \dots + Y_K(\pi)]$$

$$\pi^* = \operatorname{argmax}_{\pi} V^\pi$$

Optimal policy

Optimal treatment policy: value function

$Y_k(\pi)$: potential outcome at stage k had policy π been followed

Value function of π :

$$V^\pi = E[Y_1(\pi) \dots + Y_K(\pi)]$$

Not observed
random
variables

$$\pi^* = \operatorname{argmax}_{\pi} V^\pi$$

Outline

- Example: sepsis
- Problem formulation

A. Mathematical formulation

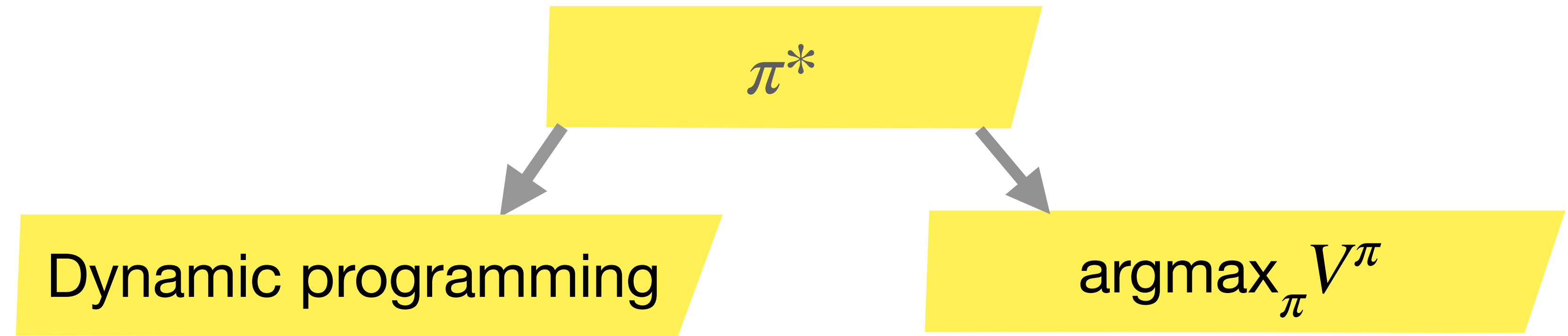
B. Existing approaches

- Proposed method
- Open questions

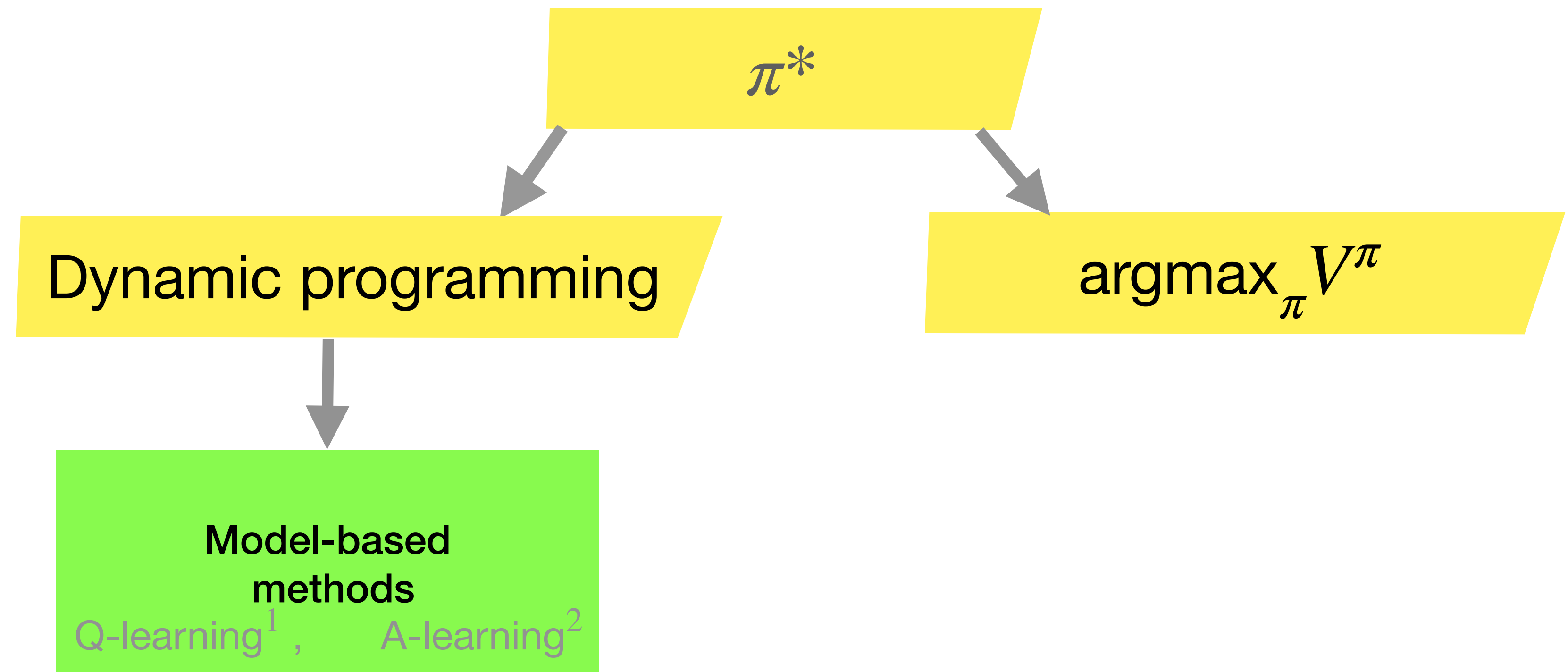
Existing approaches

Estimation of π^*

Estimation of π^*

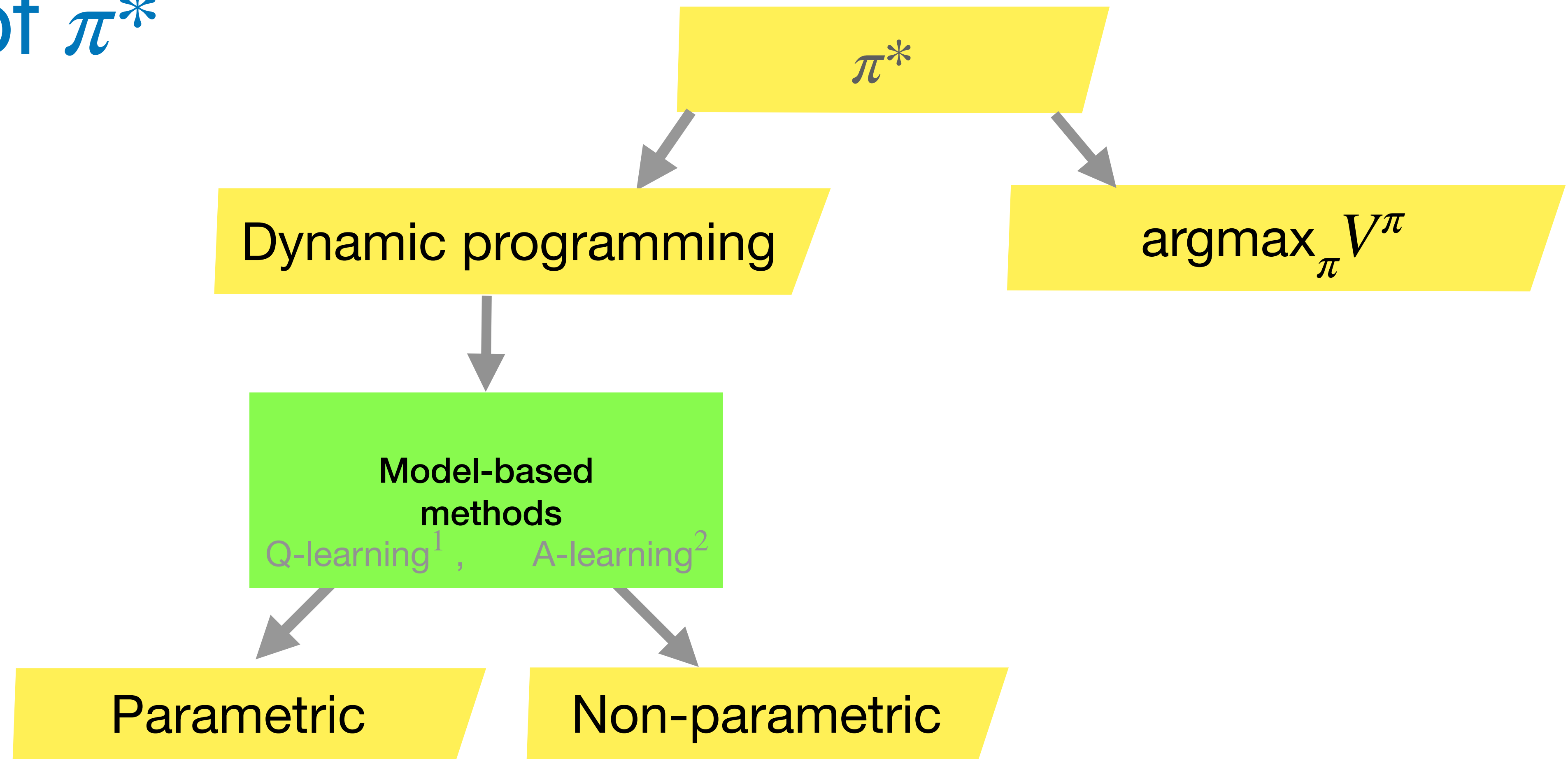


Estimation of π^*



1. Watkins, 1989; Schulte et al. 2014
2. Murphy, 2003; Robins, 2004

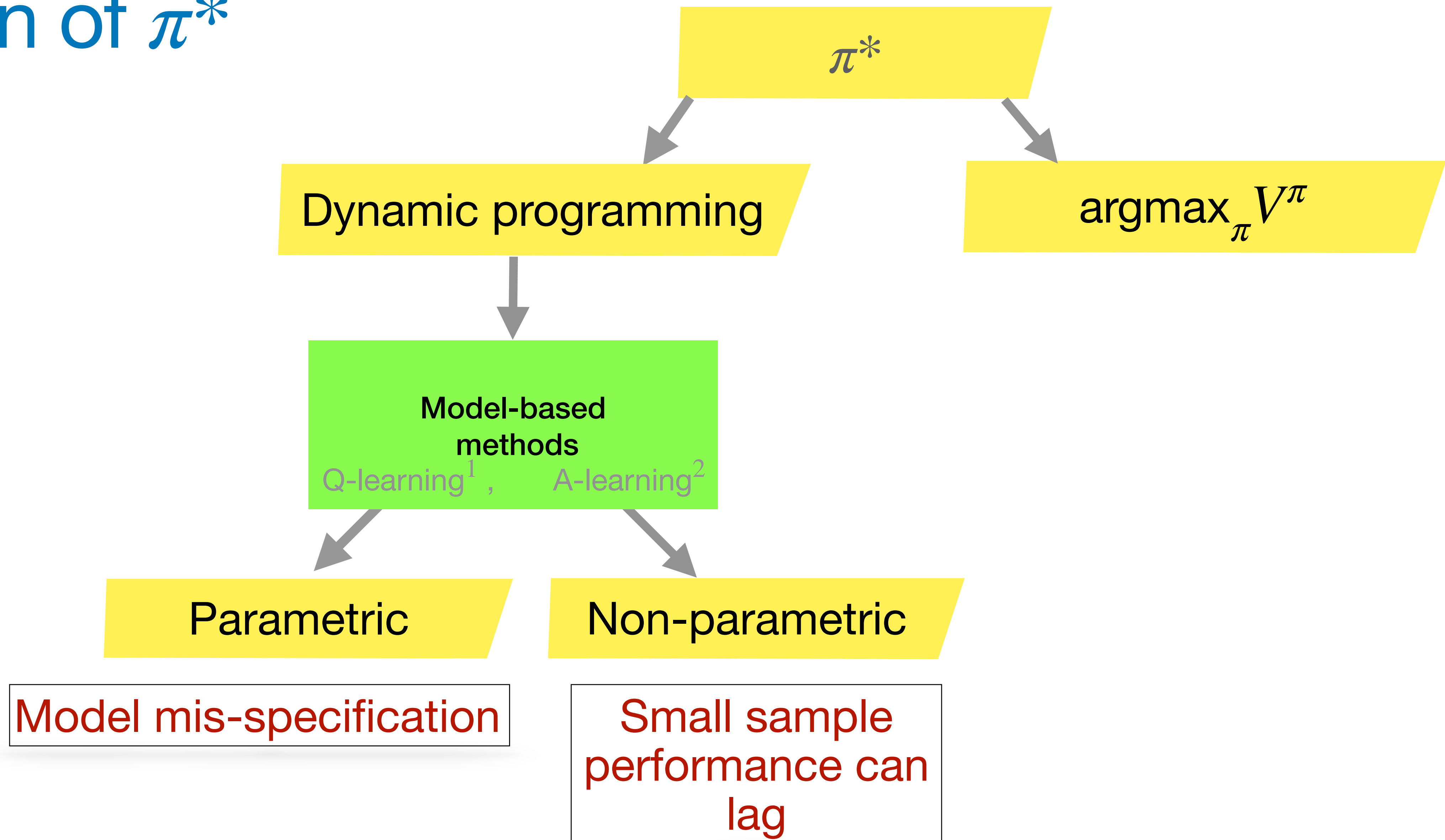
Estimation of π^*



1. Watkins, 1989; Schulte et al. 2014

2. Murphy, 2003; Robins, 2004

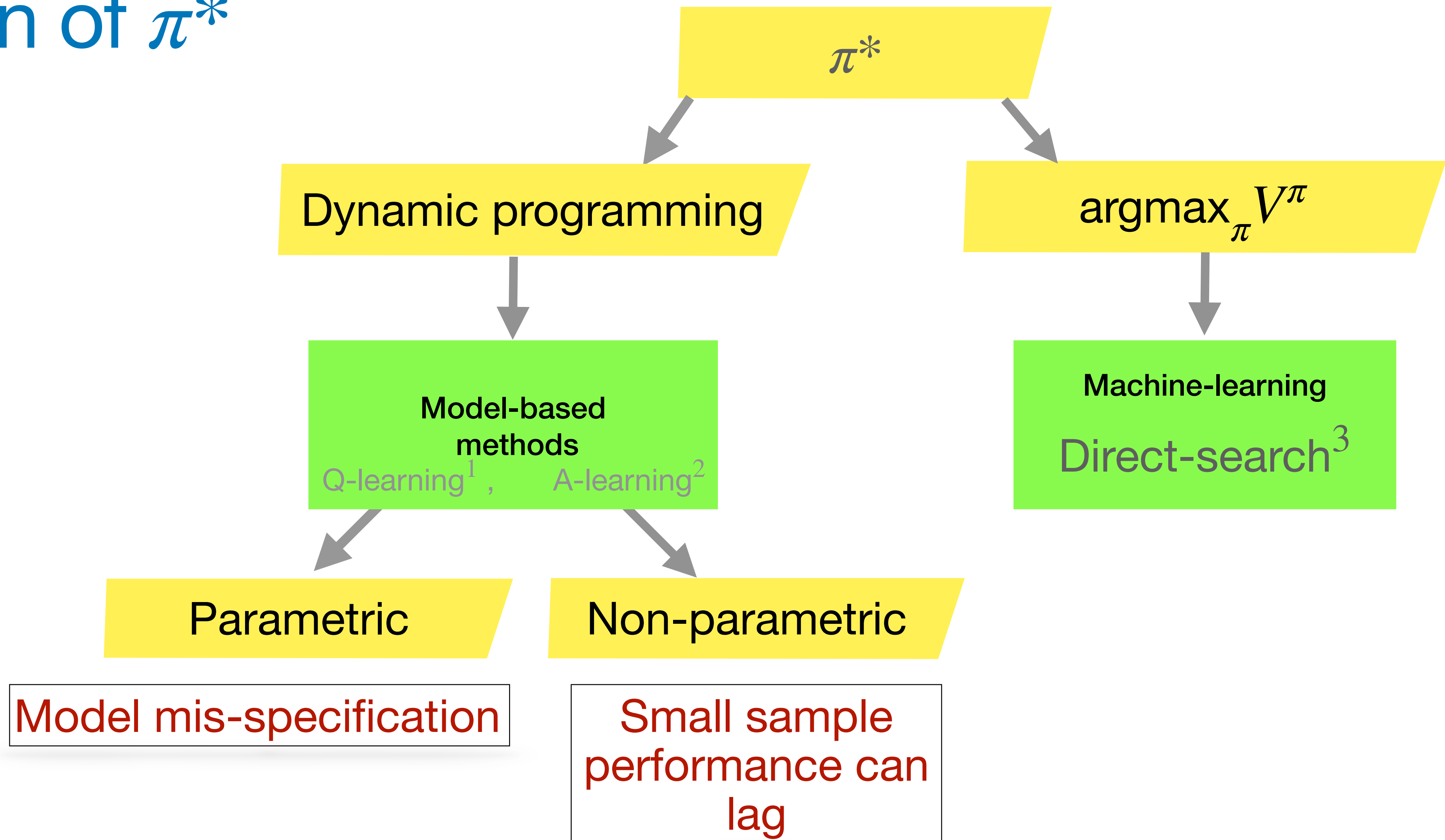
Estimation of π^*



1. Watkins, 1989; Schulte et al. 2014

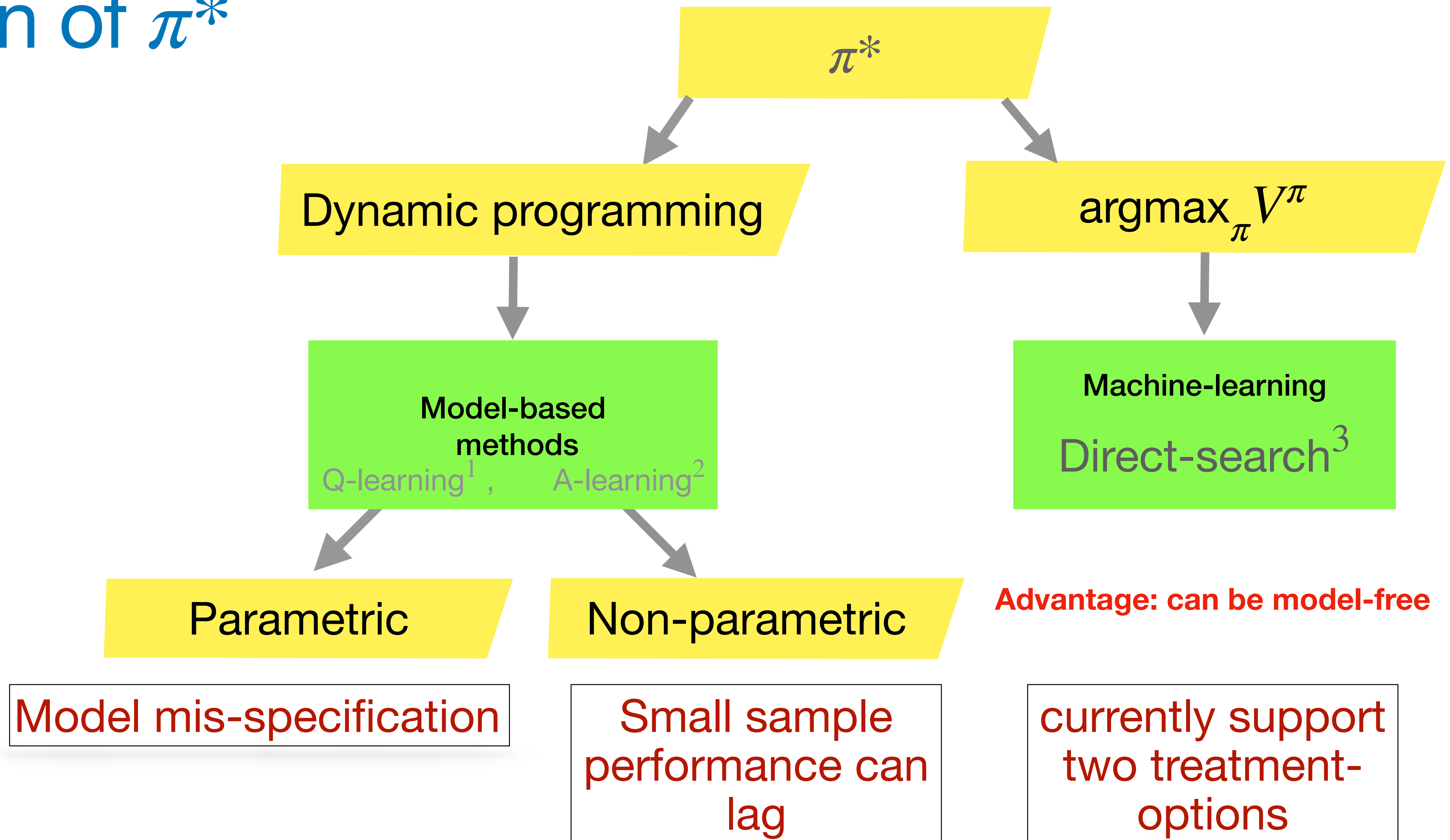
2. Murphy, 2003; Robins, 2004

Estimation of π^*



1. Watkins, 1989; Schulte et al. 2014
2. Murphy, 2003; Robins, 2004
3. Zhao et al. 2012; 2015, Laha et al. 2022

Estimation of π^*



1. Watkins, 1989; Schulte et al. 2014
2. Murphy, 2003; Robins, 2004
3. Zhao et al. 2012; 2015, Laha et al. 2022

Estimation of π^*

Goal of the project:

Dynamic programming

π^*

$\operatorname{argmax}_{\pi} V^{\pi}$

Model-based
methods

Q-learning¹, A-learning²

Machine-learning

Direct-search³

Parametric

Non-parametric

Advantage: can be model-free

Model mis-specification

Small sample
performance can
lag

currently support
two treatment-
options

1. Watkins, 1989; Schulte et al. 2014
2. Murphy, 2003; Robins, 2004
3. Zhao et al. 2012; 2015, Laha et al. 2022

Estimation of π^*

Goal of the project:

1. direct search for arbitrary number of treatments

Dynamic programming

Model-based methods

Q-learning¹, A-learning²

Parametric

Model mis-specification

Non-parametric

Small sample performance can lag

π^*

$\operatorname{argmax}_{\pi} V^{\pi}$

Machine-learning

Direct-search³

Advantage: can be model-free

currently support two treatment-options

1. Watkins, 1989; Schulte et al. 2014
2. Murphy, 2003; Robins, 2004
3. Zhao et al. 2012; 2015, Laha et al. 2022

Estimation of π^*

Goal of the project:

1. direct search for arbitrary number of treatments
2. Computationally efficient and scalable

Dynamic programming

Model-based methods

Q-learning¹, A-learning²

Parametric

Model mis-specification

Non-parametric

Small sample performance can lag

π^*

$\operatorname{argmax}_{\pi} V^{\pi}$

Machine-learning

Direct-search³

Advantage: can be model-free

currently support two treatment-options

1. Watkins, 1989; Schulte et al. 2014
2. Murphy, 2003; Robins, 2004
3. Zhao et al. 2012; 2015, Laha et al. 2022

Outline

- Example: sepsis
- Problem formulation
- Proposed method

A. Methodology

B. Example on a toy data

- Open questions

Outline

- Example: sepsis
- Problem formulation
- Proposed method

A. Methodology

B. Example on a toy data

- Open questions

Proposed method

Only two treatment options

Only two treatment options

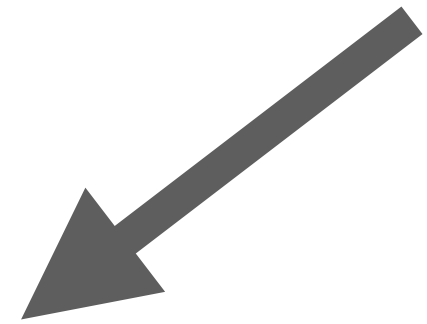


H_1

Only two treatment options



H_1



Treatment



Only two treatment options



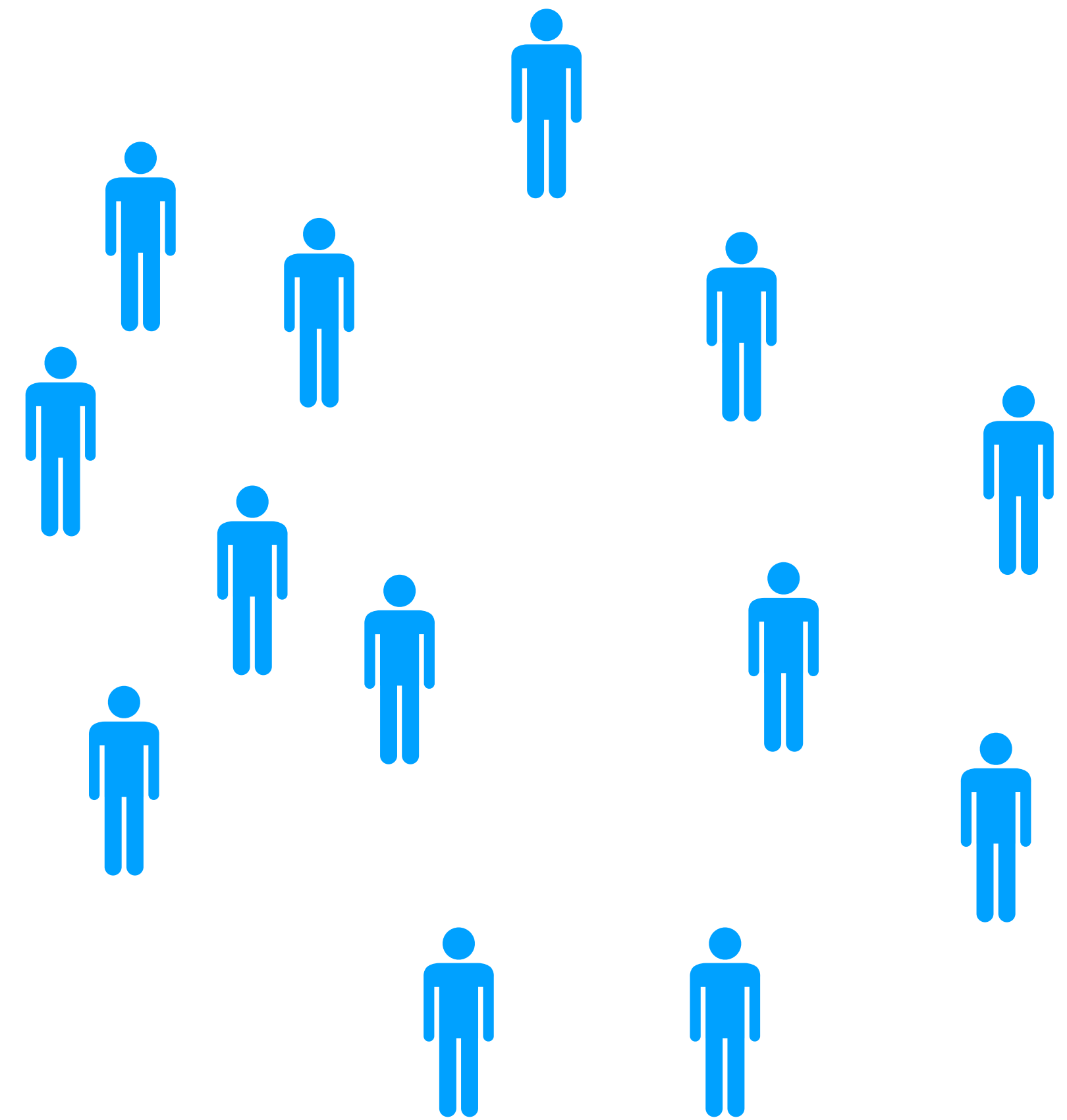
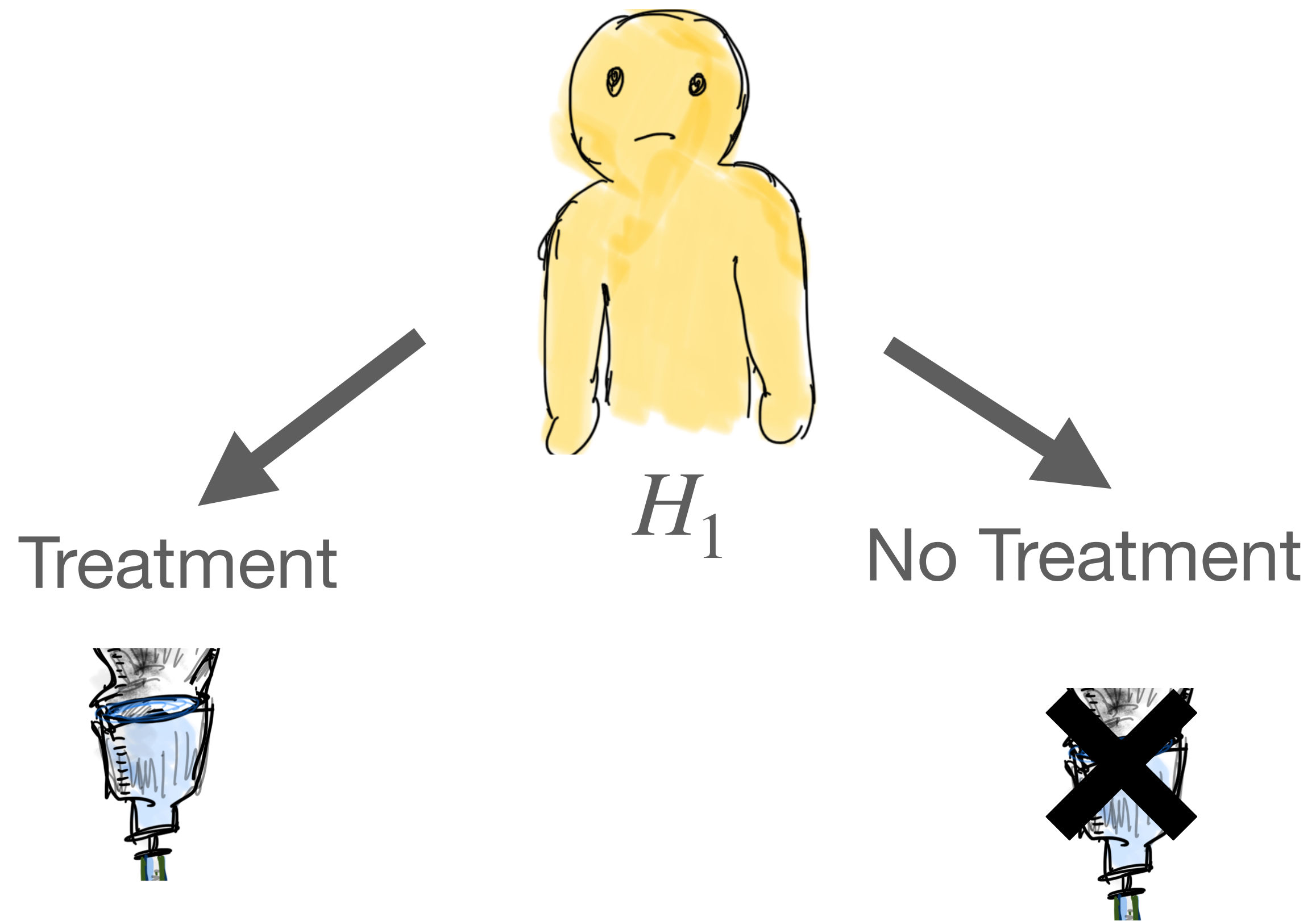
H_1

Treatment

No Treatment

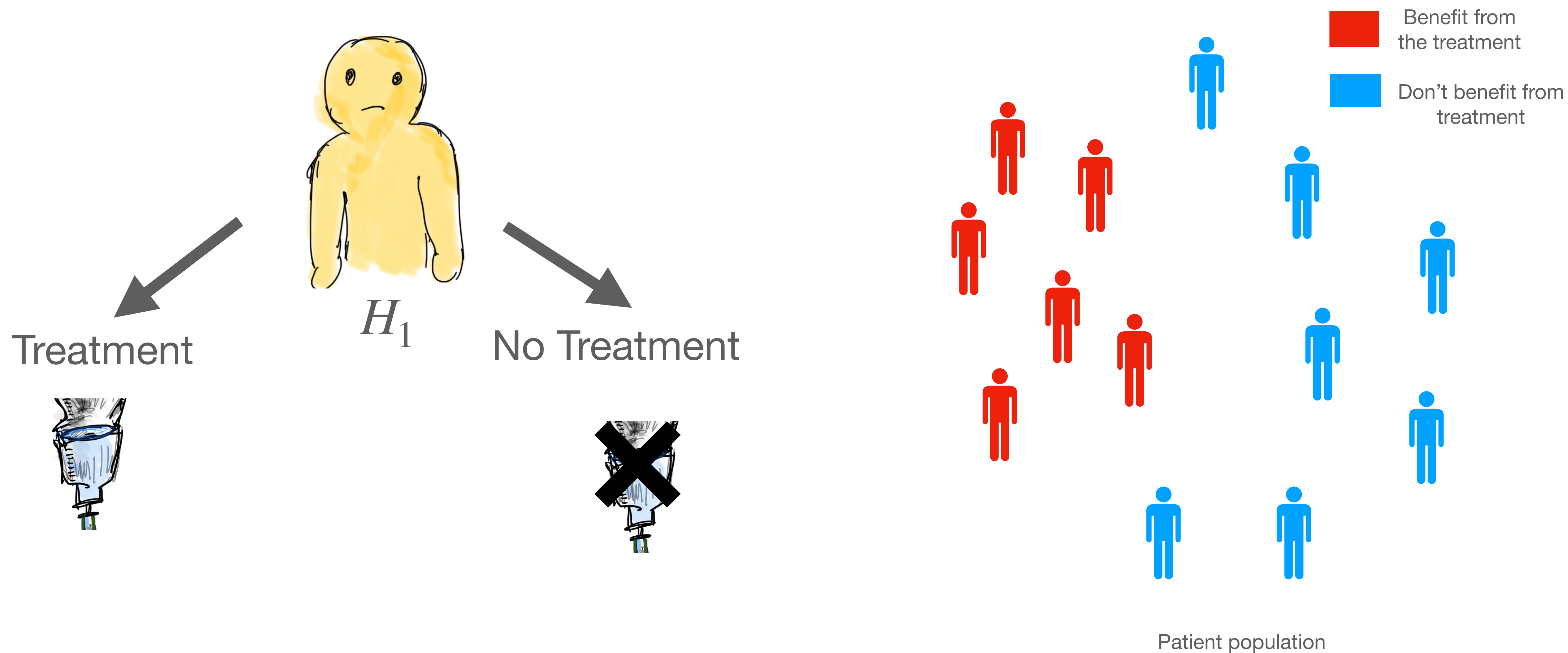


Only two treatment options



Patient population

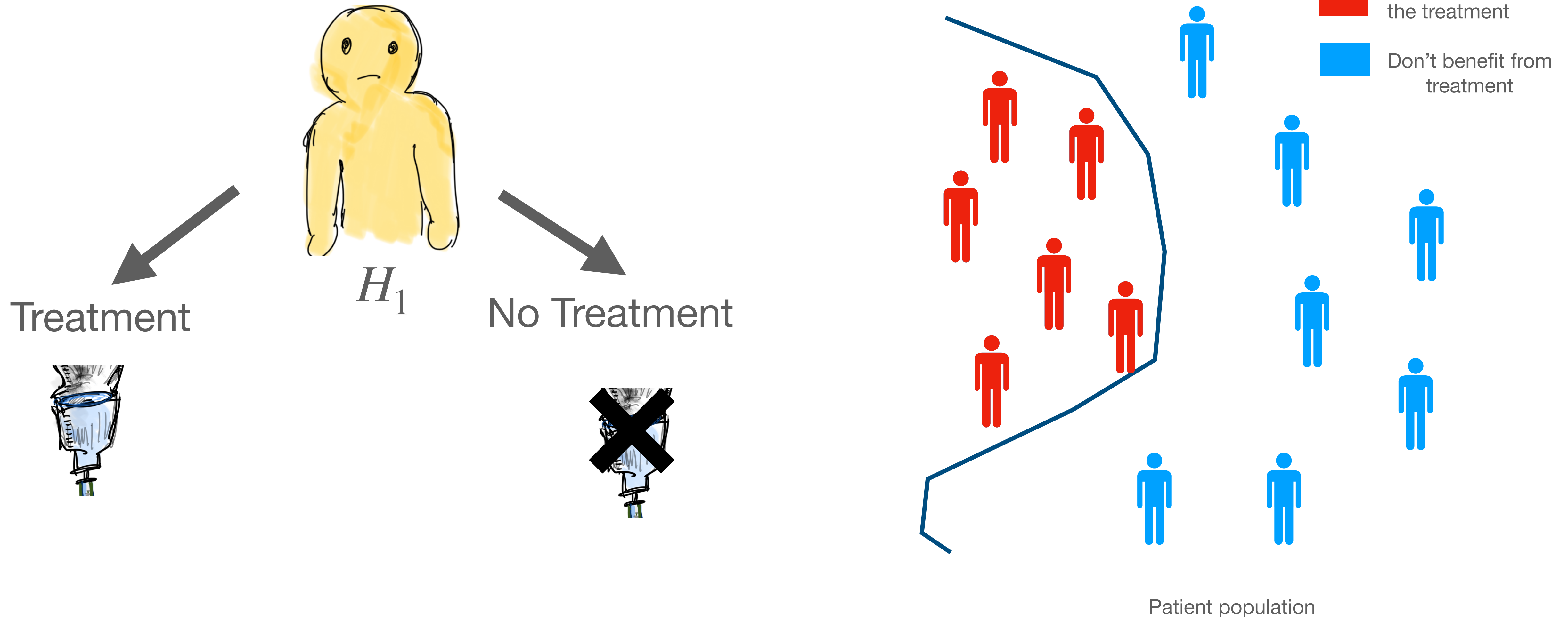
Only two treatment options



Only two treatment options

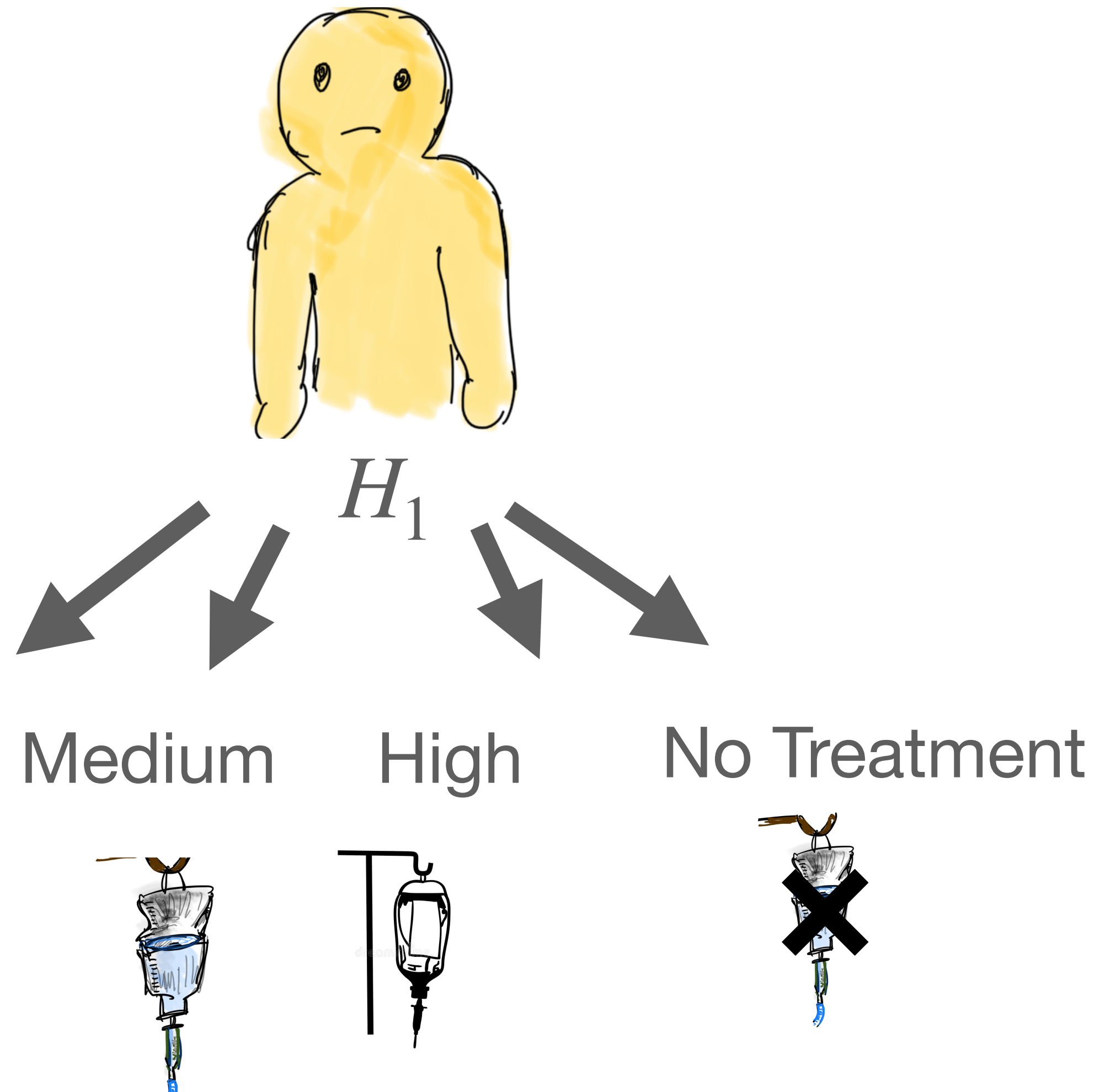
Treatment assignment at each stage:

binary classification problem

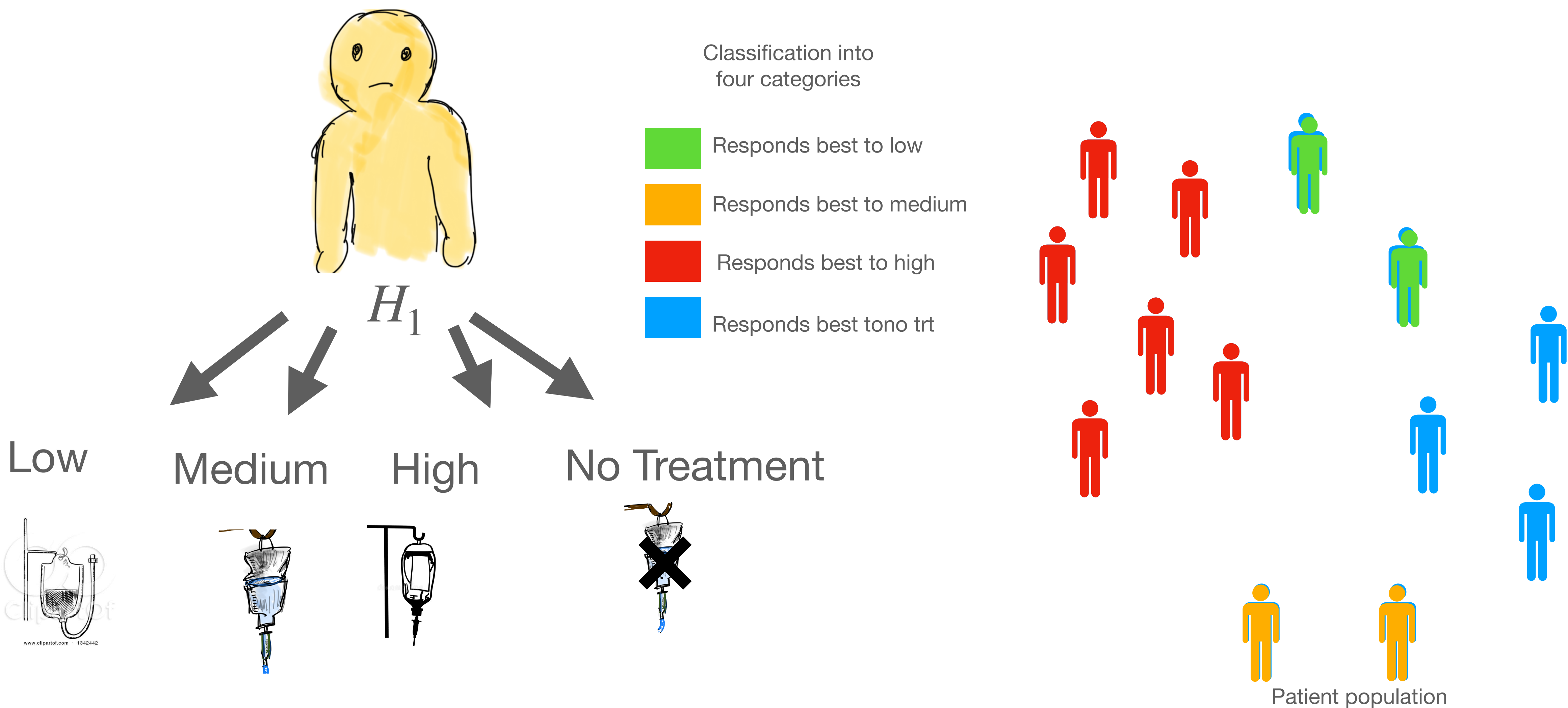


More than two treatment option

More than two treatment option

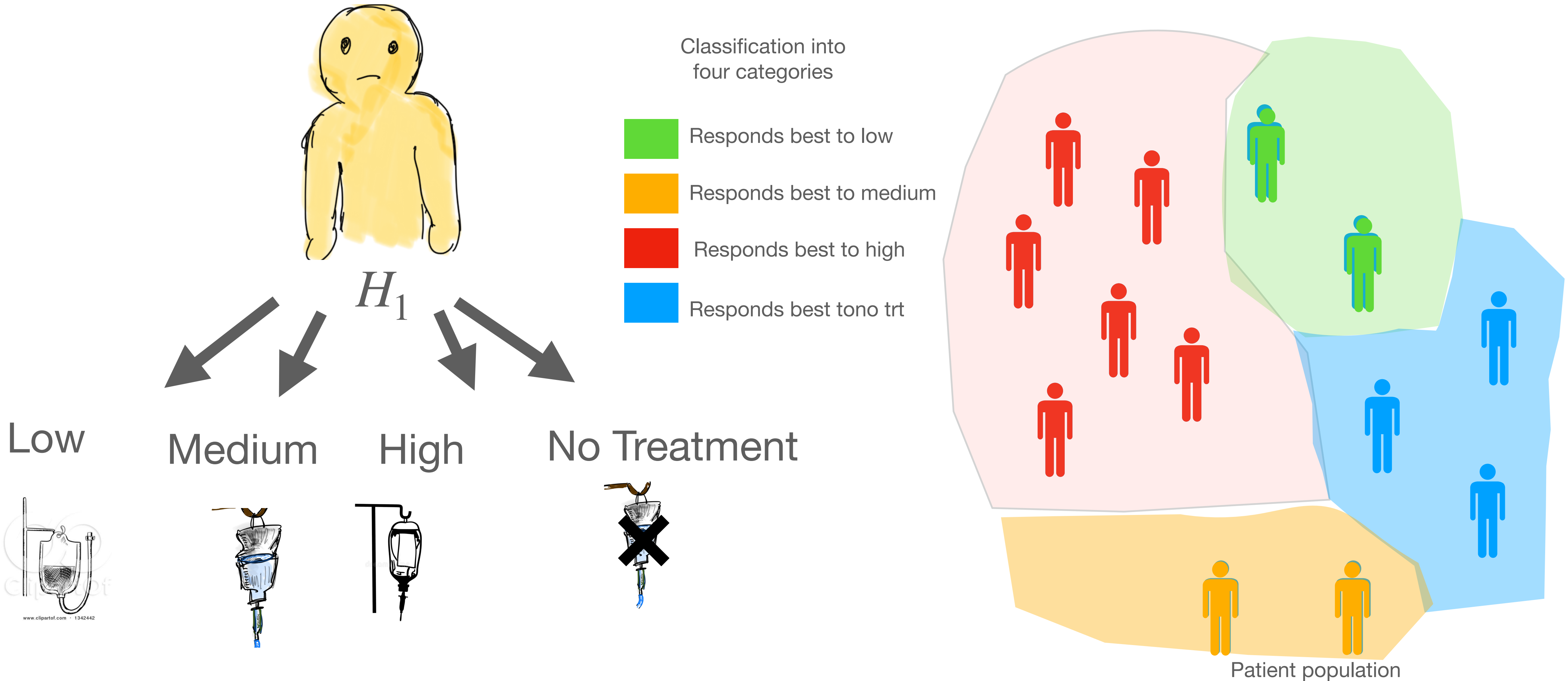


More than two treatment option



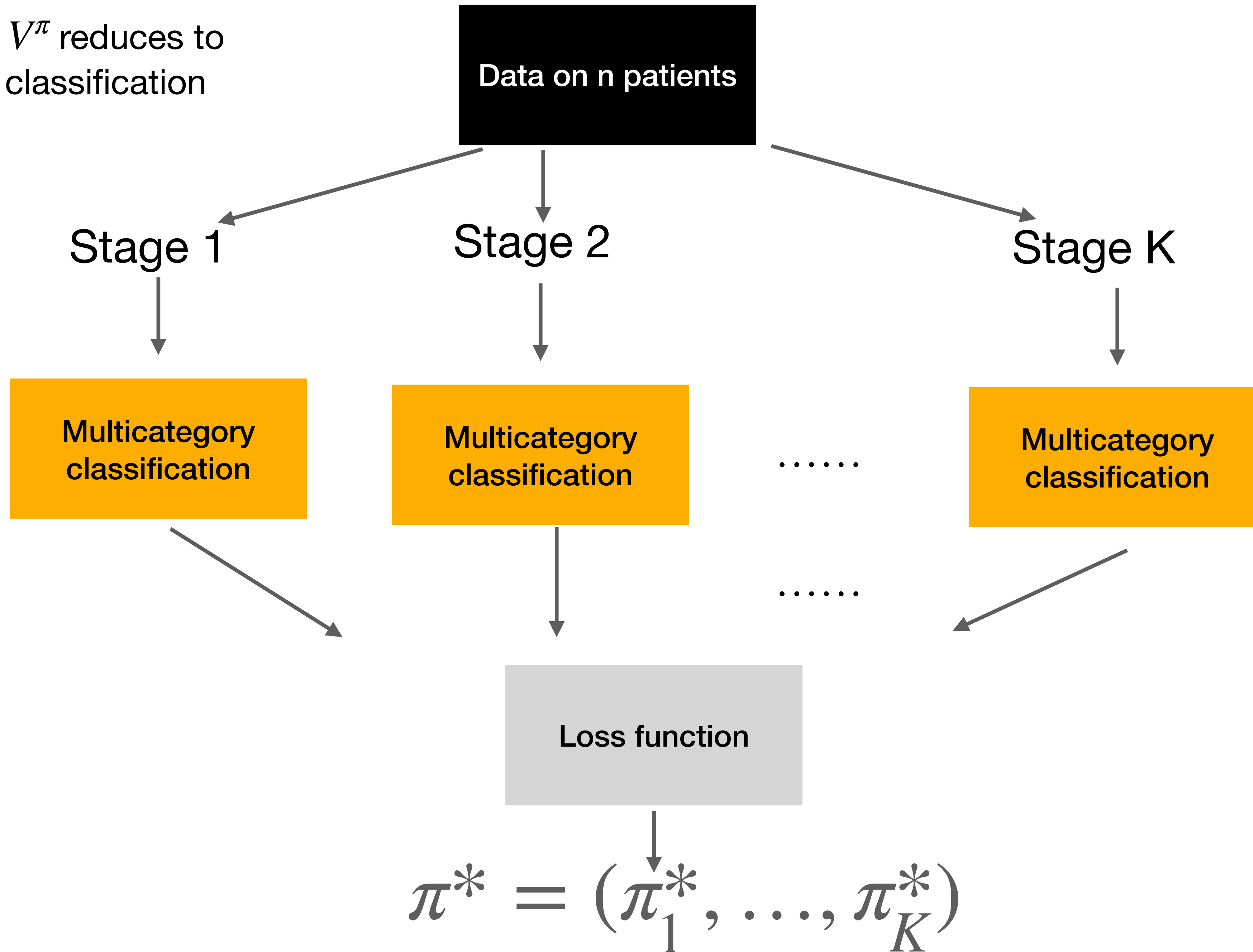
More than two treatment option

Treatment assignment at each stage: multcategory classification problem



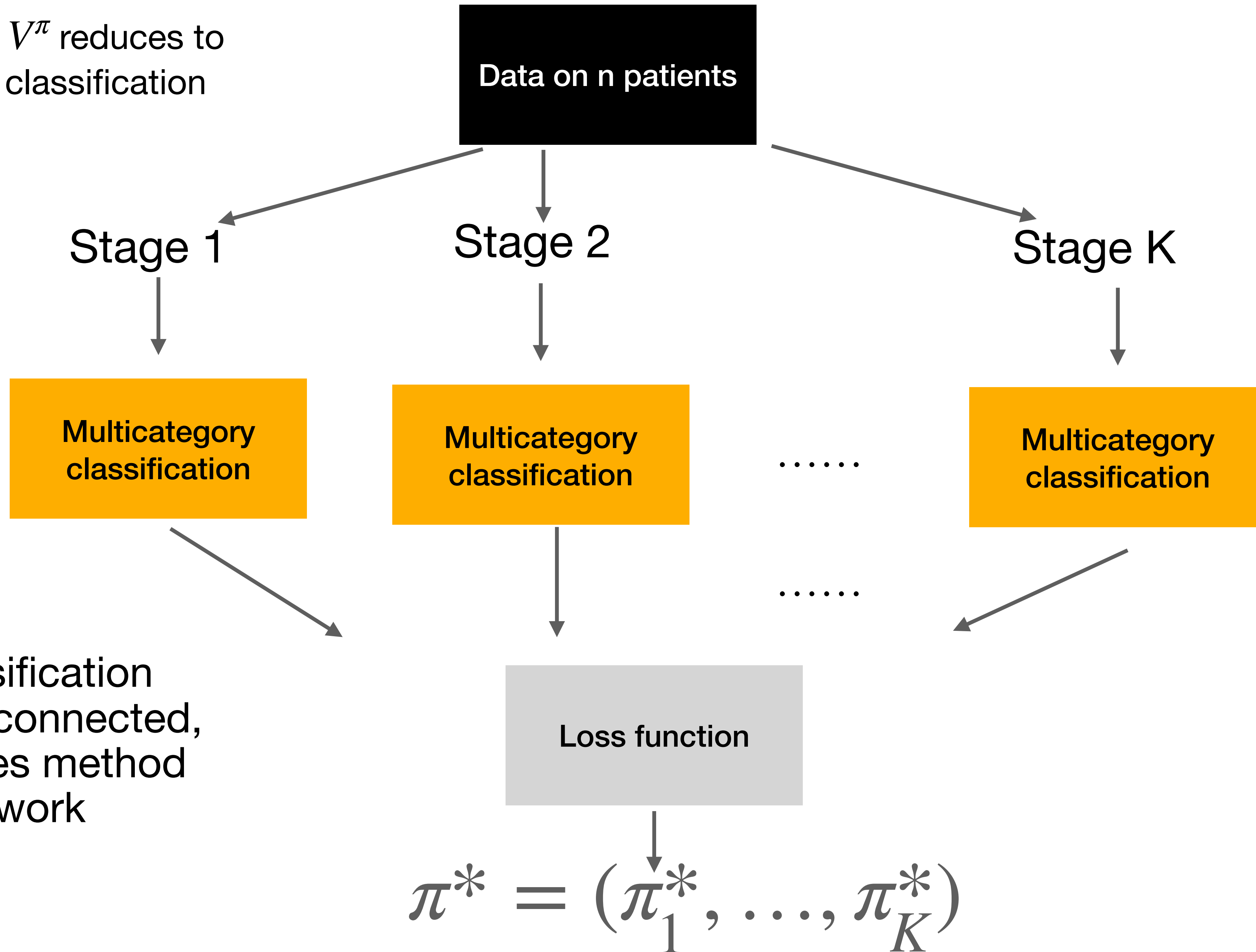
Proposed method

Maximization of V^π reduces to simultaneous K classification problems



Proposed method

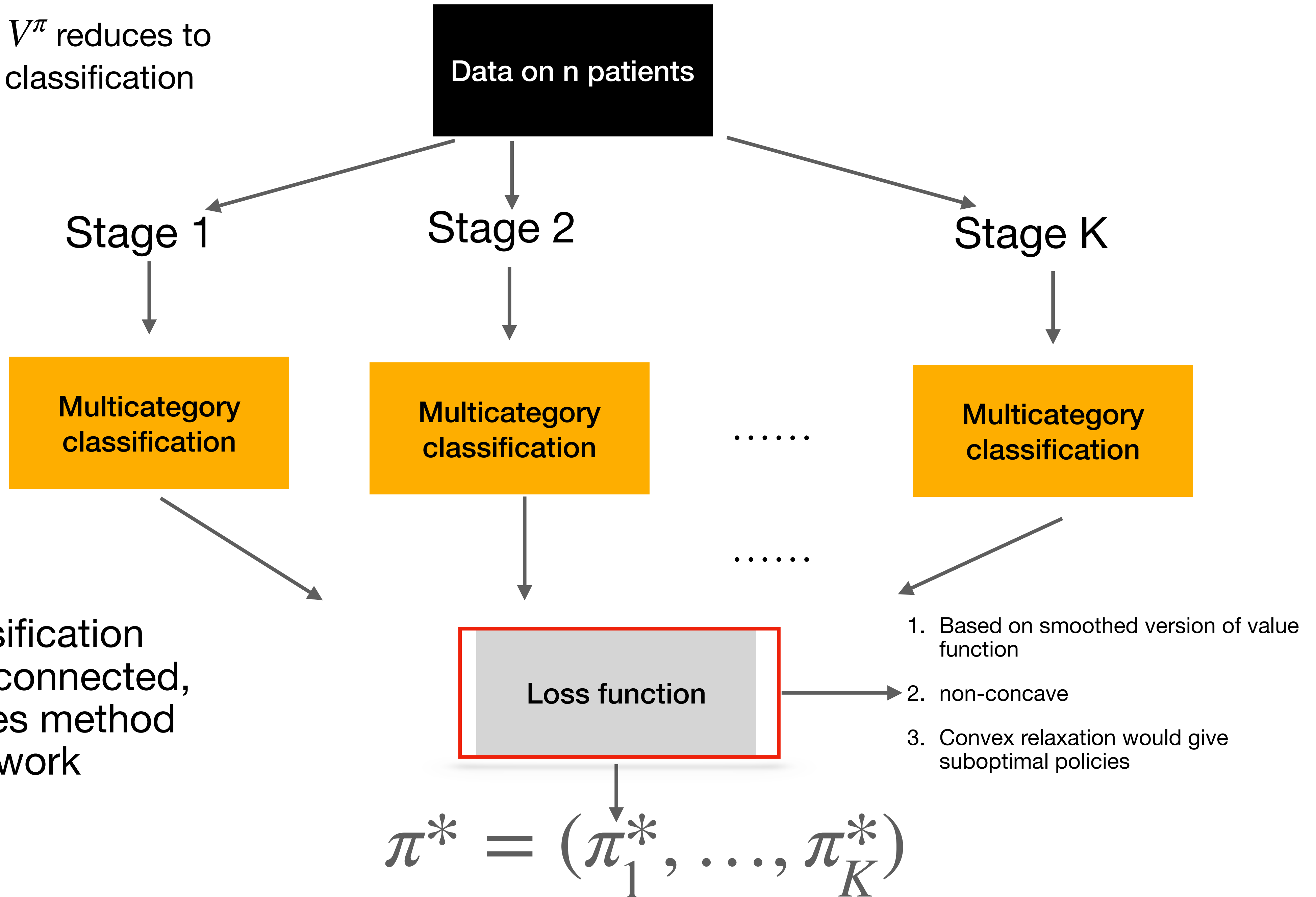
Maximization of V^π reduces to simultaneous K classification problems



The K classification problems are connected, off-the-shelves method will not work

Proposed method

Maximization of V^π reduces to simultaneous K classification problems

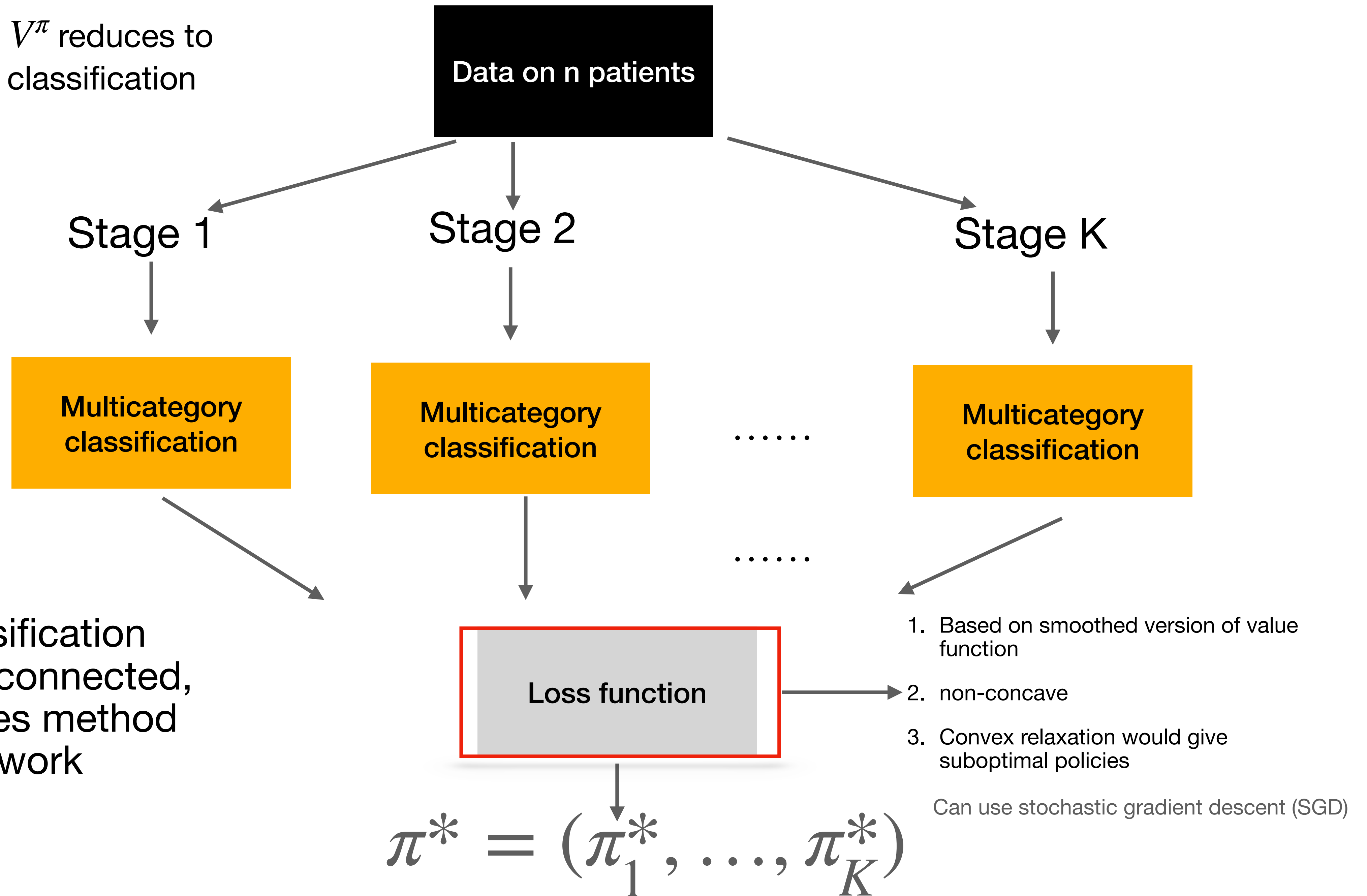


The K classification problems are connected, off-the-shelves method will not work

$$\pi^* = (\pi_1^*, \dots, \pi_K^*)$$

Proposed method

Maximization of V^π reduces to simultaneous K classification problems



The K classification problems are connected, off-the-shelves method will not work

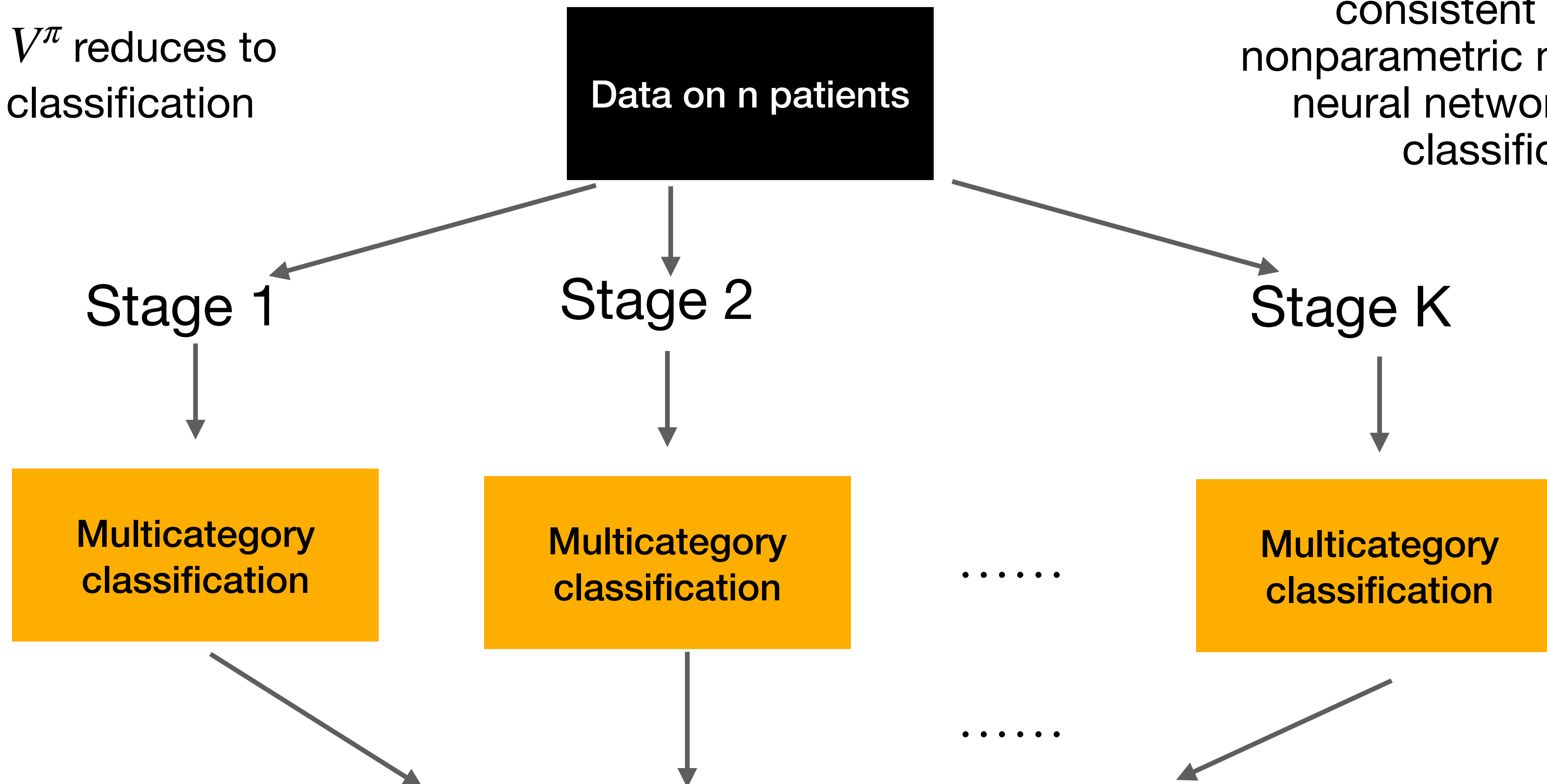
$$\pi^* = (\pi_1^*, \dots, \pi_K^*)$$

Proposed method

Maximization of V^π reduces to simultaneous K classification problems

Estimated policy will be consistent if we use nonparametric methods, e.g., neural networks, for the classification

The K classification problems are connected, off-the-shelves method will not work



1. Based on smoothed version of value function
2. non-concave
3. Convex relaxation would give suboptimal policies

Can use stochastic gradient descent (SGD)

Population level solution

$$\pi^* = (\pi_1^*, \dots, \pi_K^*)$$

Outline

- Example: sepsis
- Problem formulation
- Proposed method

A. Methodology

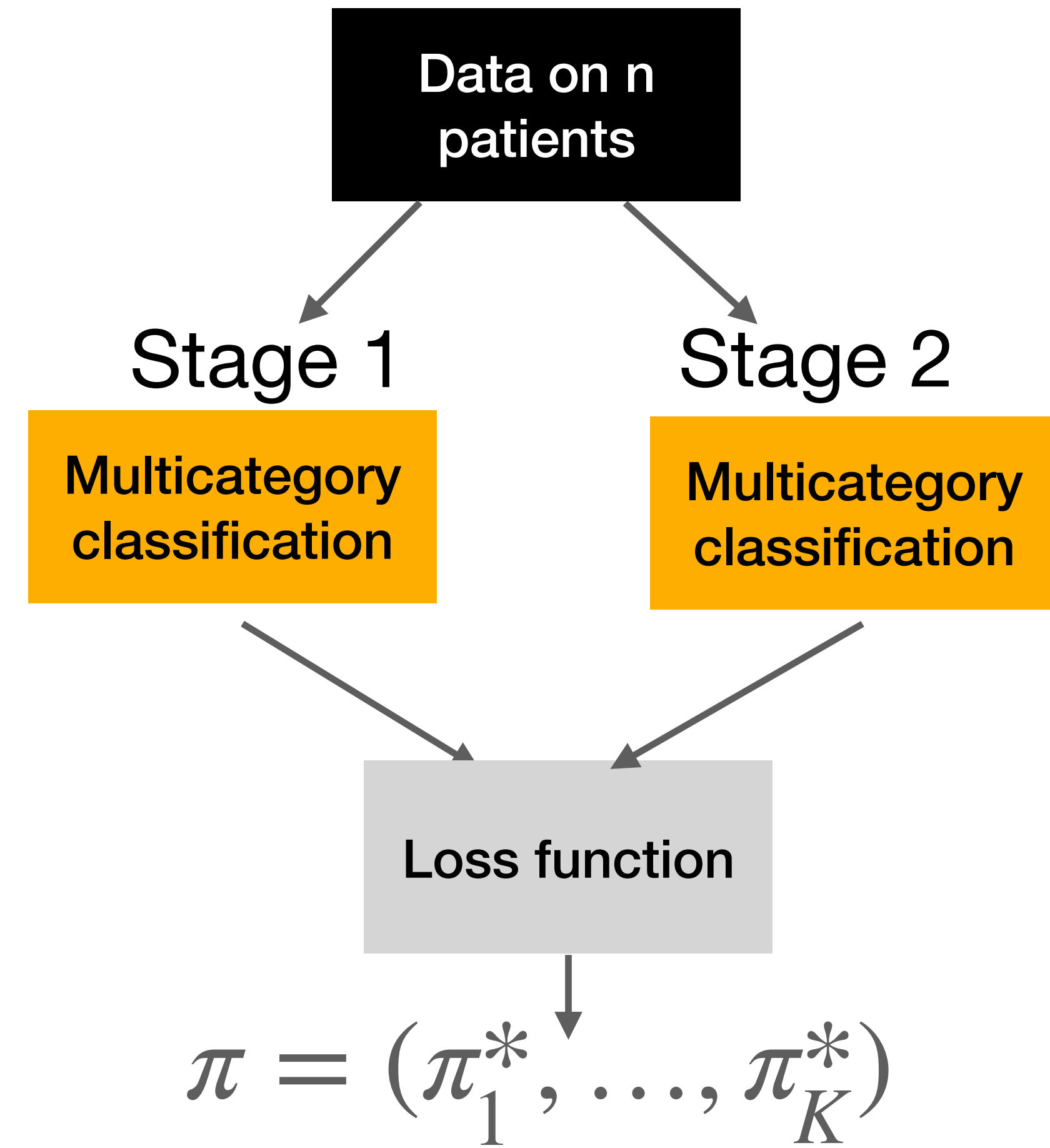
B. Example on a toy data

- Open questions

Example with toy data

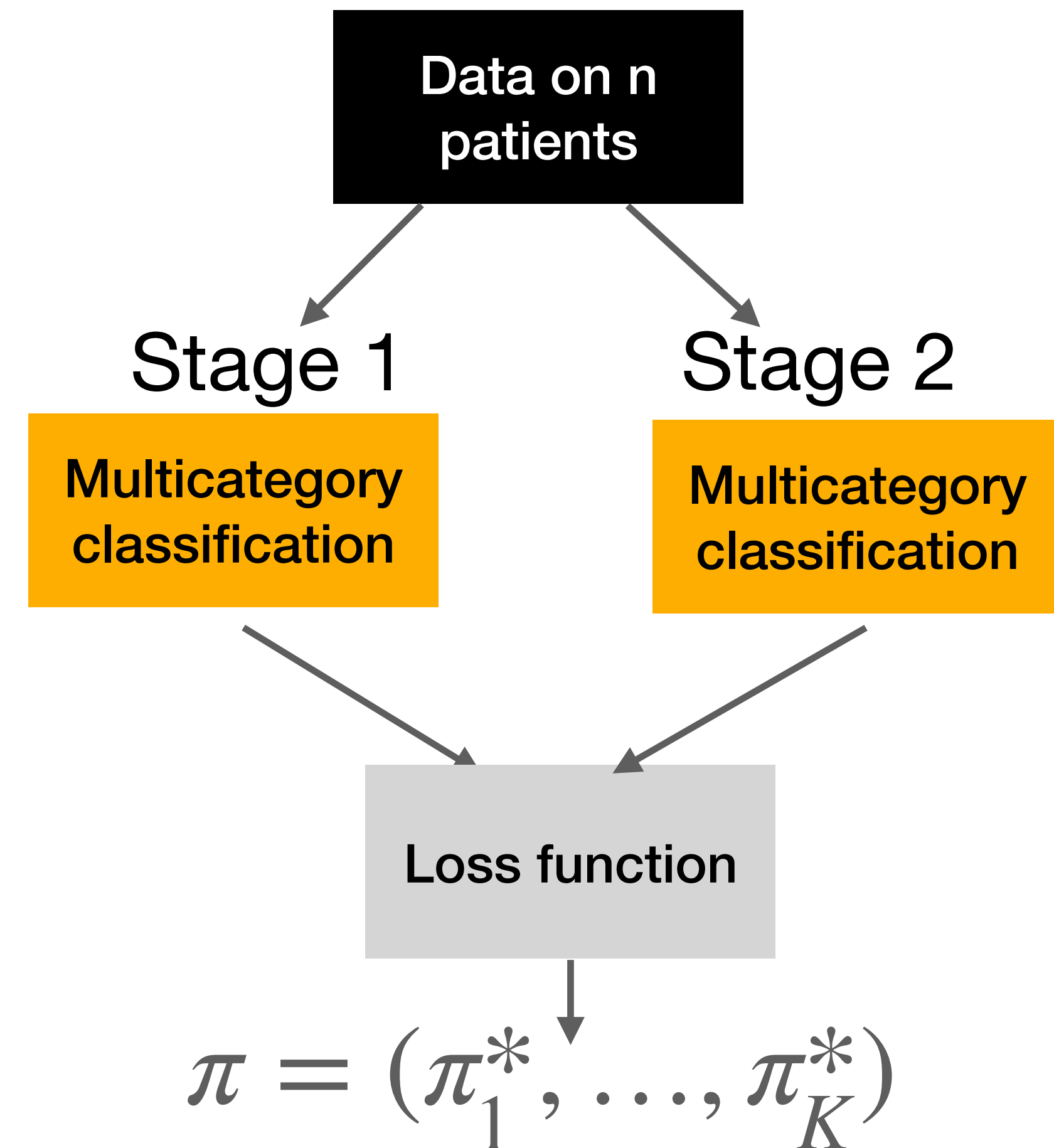
*This work is by Sneha Mishra, my former summer RA

One example



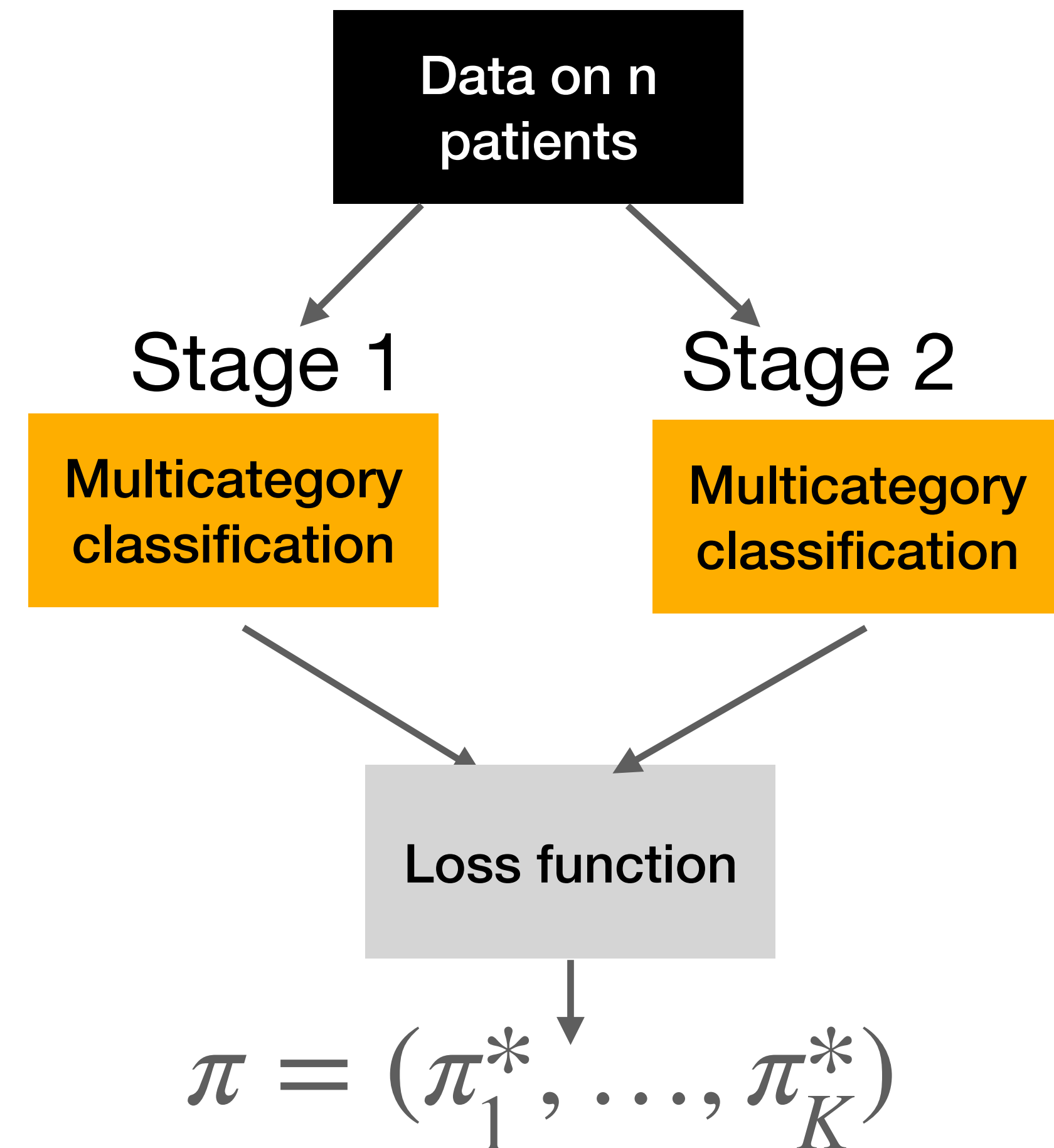
One example

- Suppose number of stages, i.e., $K = 2$



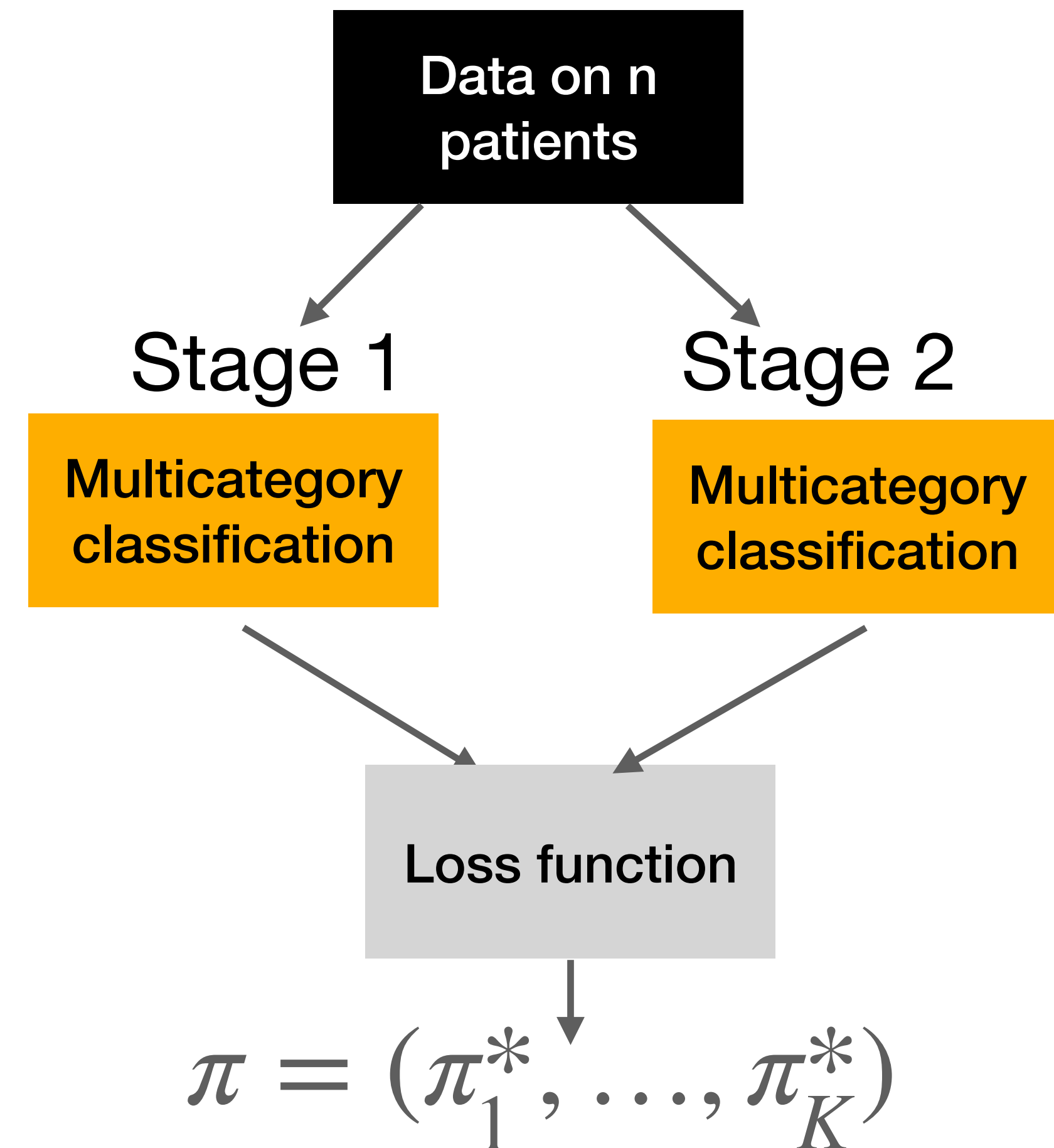
One example

- Suppose number of stages, i.e., $K = 2$
- Number of treatments at each stage: 3.



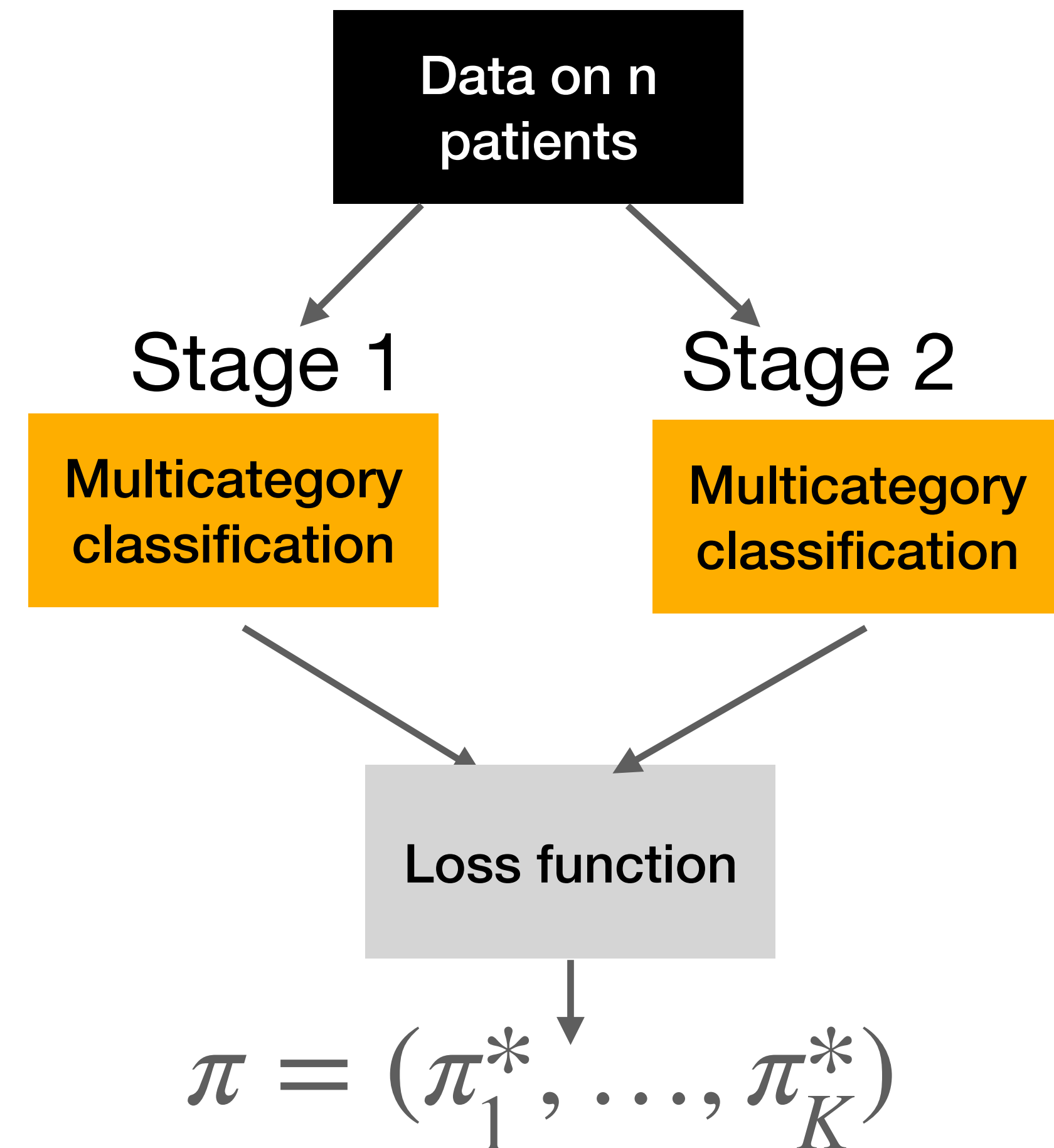
One example

- Suppose number of stages, i.e., $K = 2$
- Number of treatments at each stage: 3.
- Use neural network classifiers



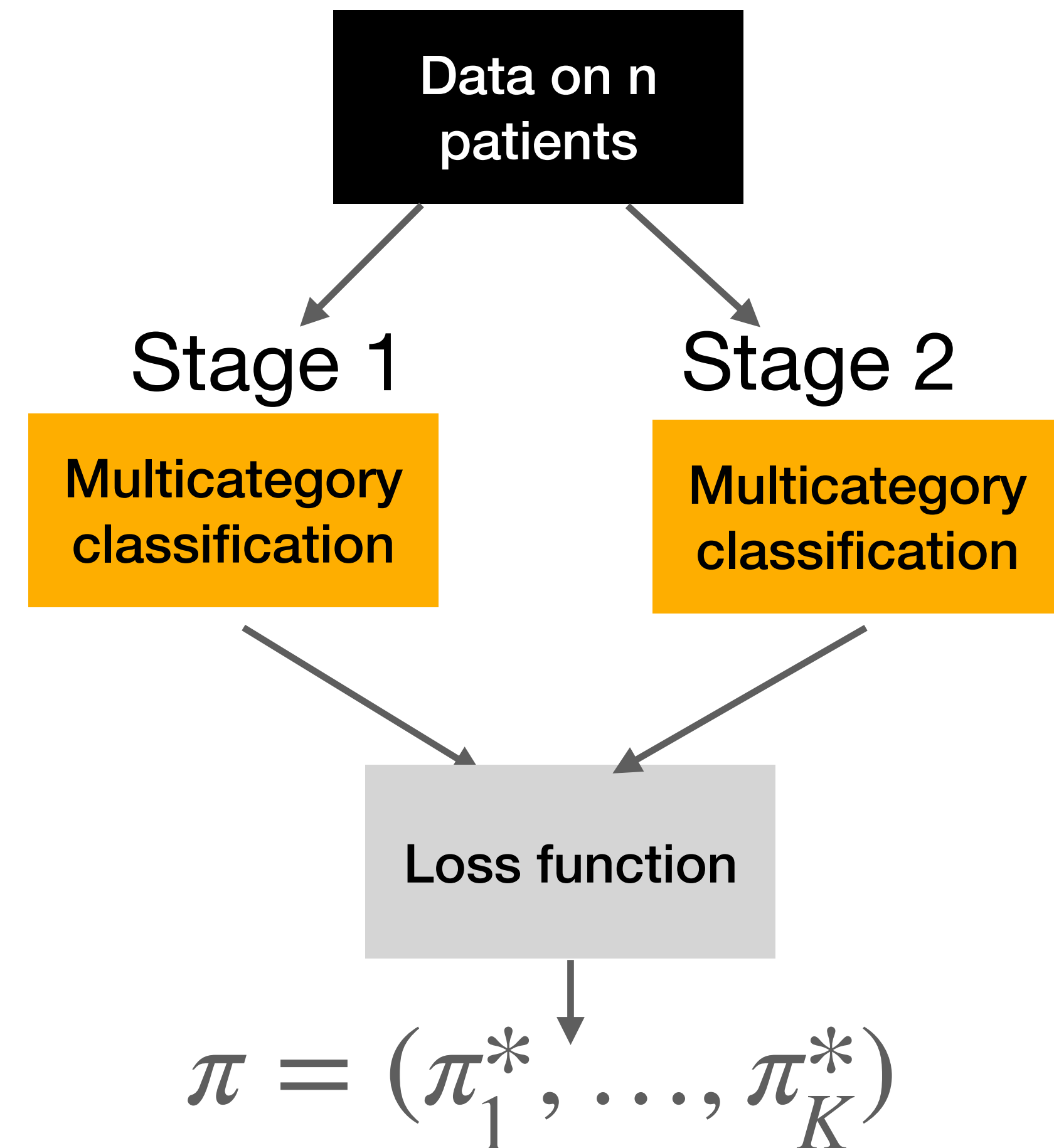
One example

- Suppose number of stages, i.e., $K = 2$
- Number of treatments at each stage: 3.
- Use neural network classifiers
- No. Of covariates: 3

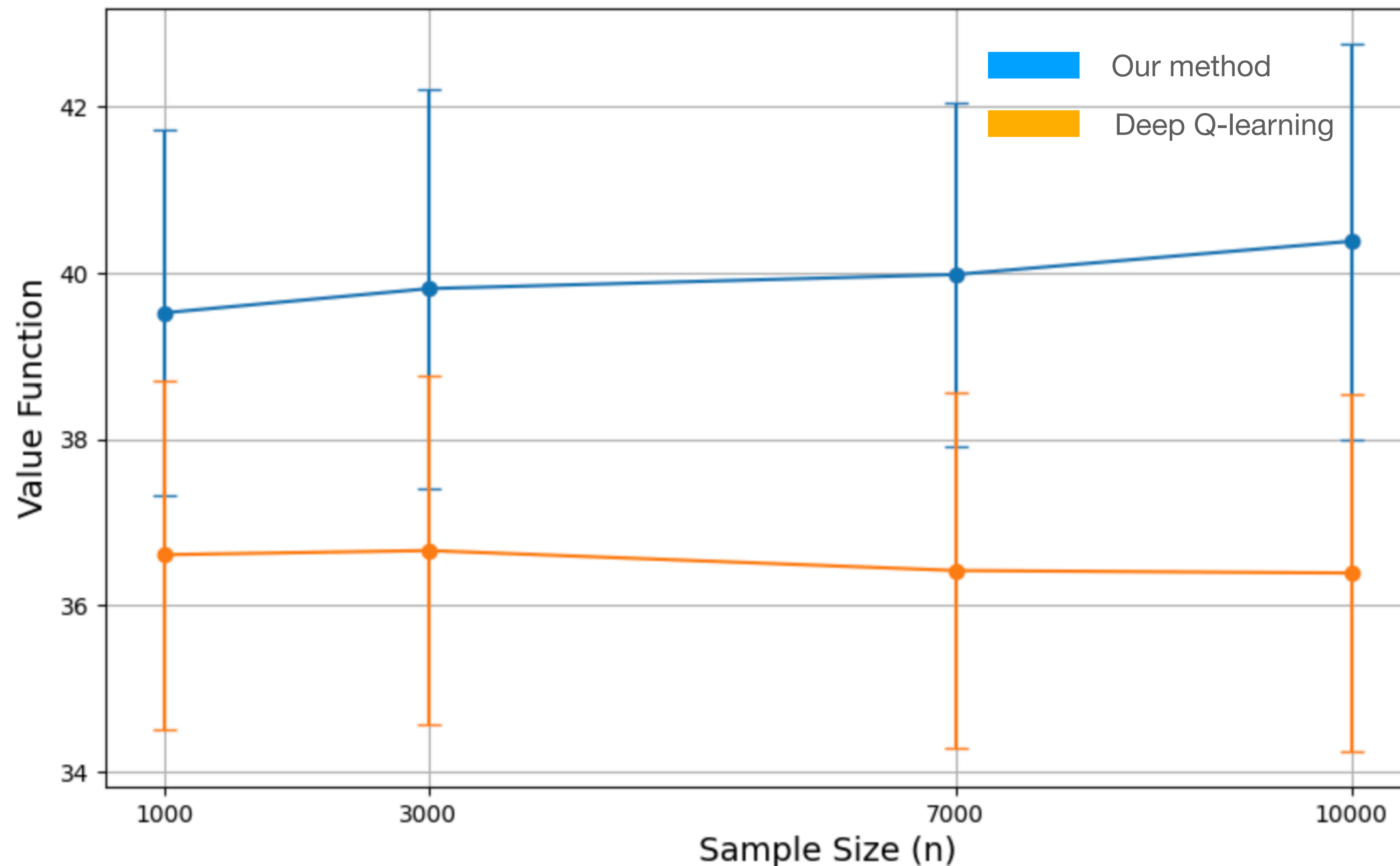


One example

- Suppose number of stages, i.e., $K = 2$
- Number of treatments at each stage: 3.
- Use neural network classifiers
- No. Of covariates: 3
- The covariates and rewards were Gaussian, and the rewards were generated by a linear model.



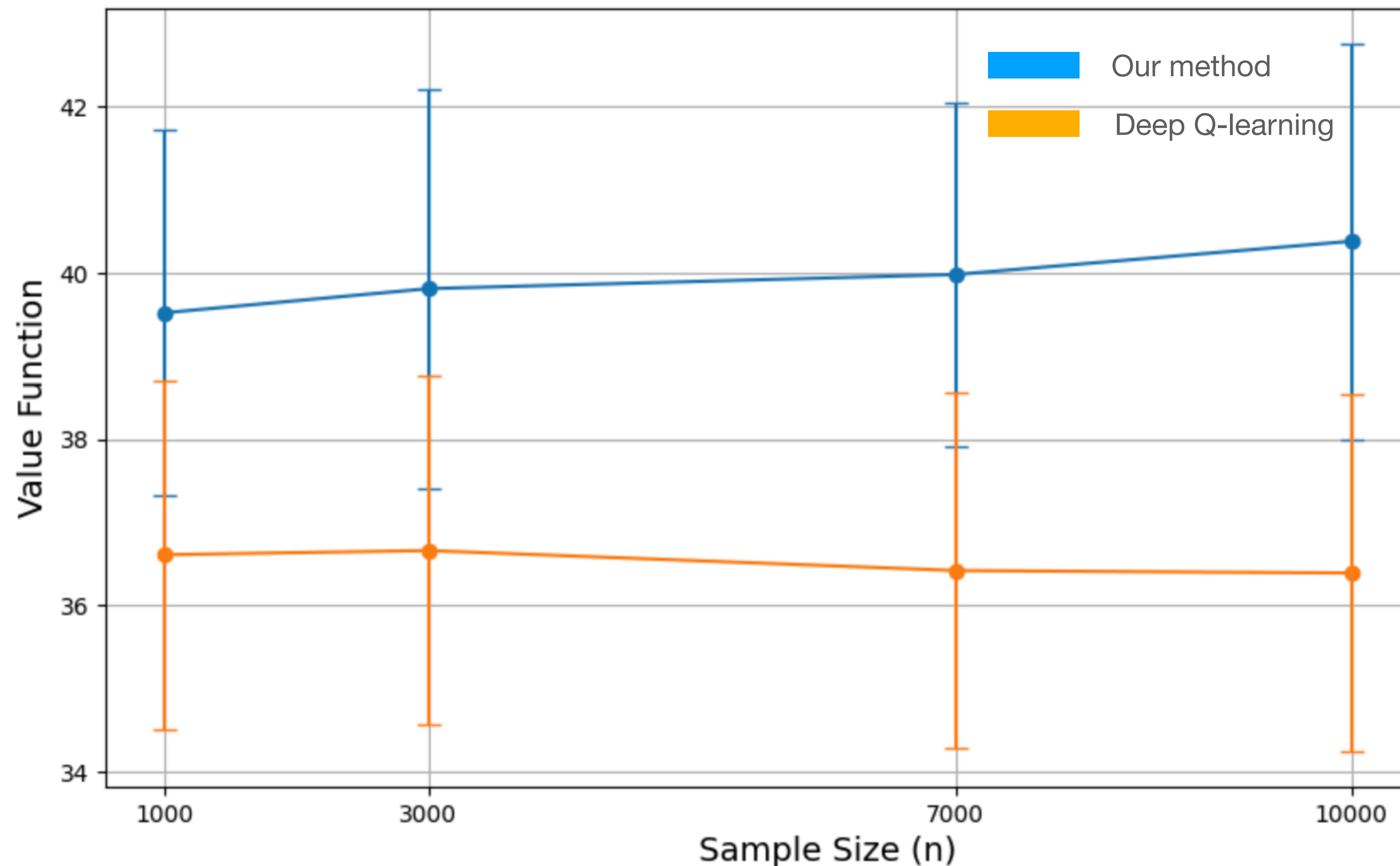
Plot of the population-level value functions



The deep Q-learning line represents the optimal policy generated by deep Q-learning method for DTR — that is current gold standard

Plot of the value function for our method and deep Q-learning based on the toy data

Plot of the population-level value functions



The deep Q-learning line represents the optimal policy generated by deep Q-learning method for DTR — that is current gold standard

Plot of the value function for our method and deep Q-learning based on the toy data

Outline

- Example: sepsis
- Problem formulation
- Proposed method

- Open questions

A. Implementation and optimization

B. Regret decay rate

C. Doubly robust learning

Outline

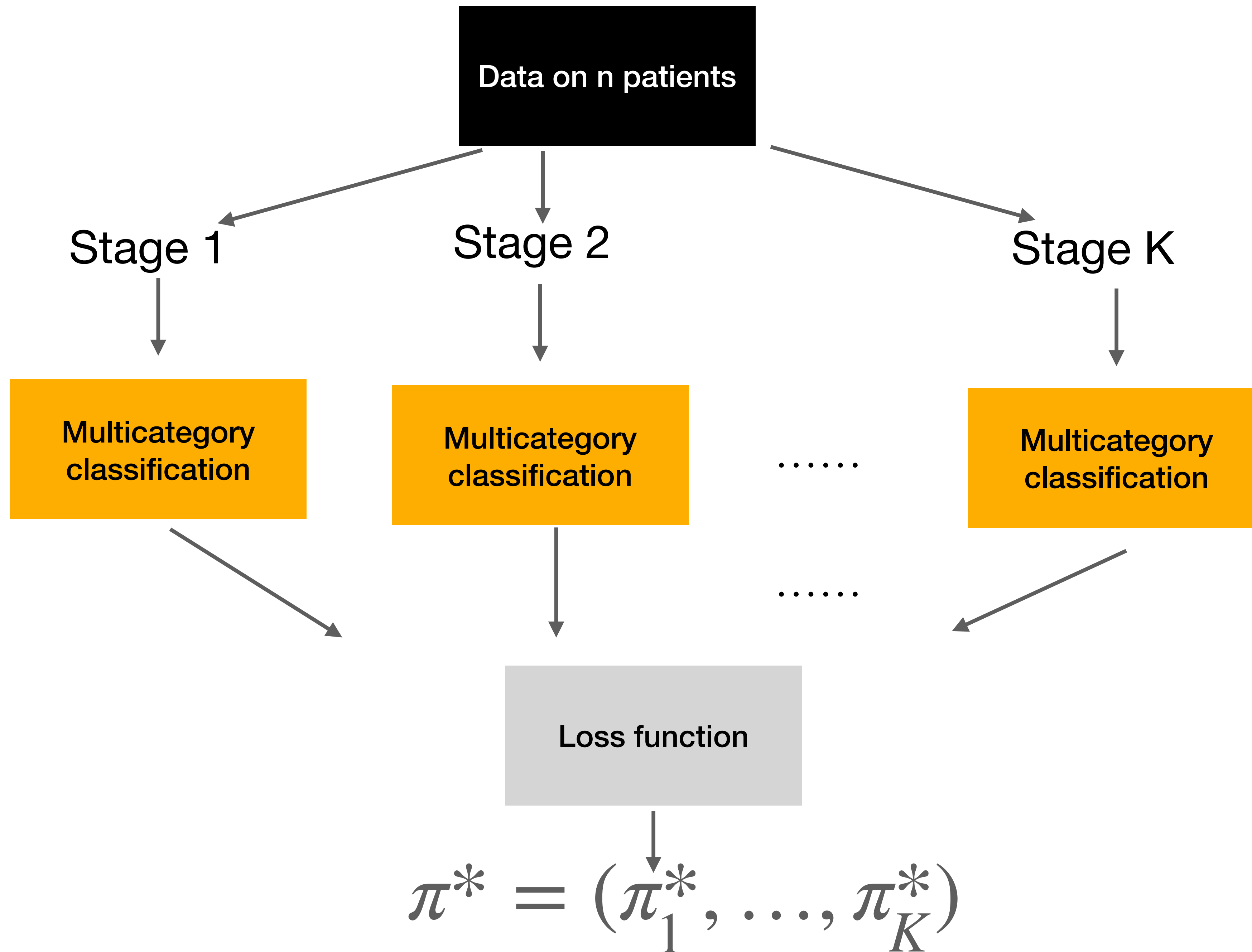
- Example: sepsis
- Problem formulation
- Proposed method
- Open questions

A. Implementation and optimization

B. Regret decay rate

C. Doubly robust learning

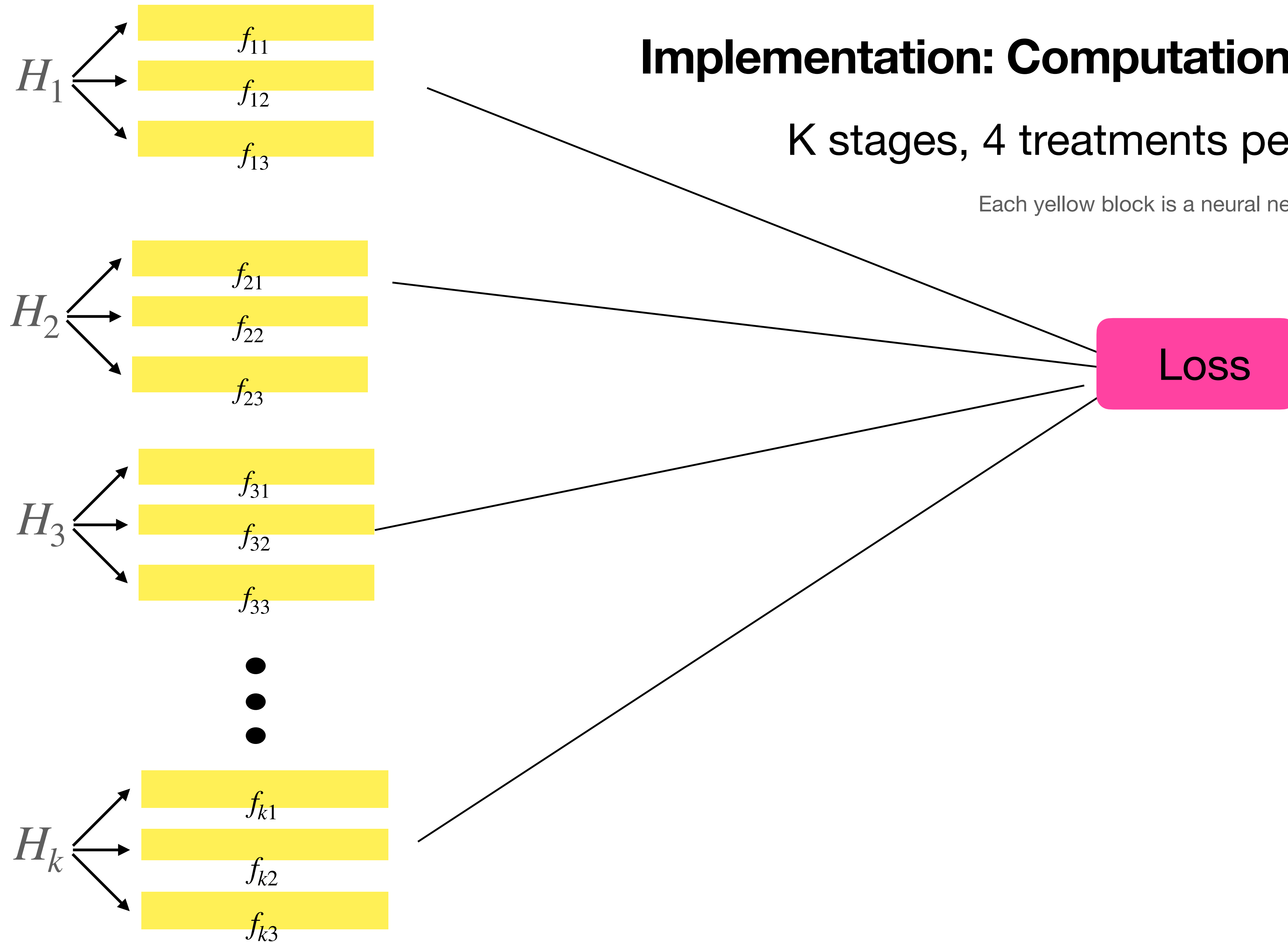
Implementation



Implementation: Computational challenge

K stages, 4 treatments per stage

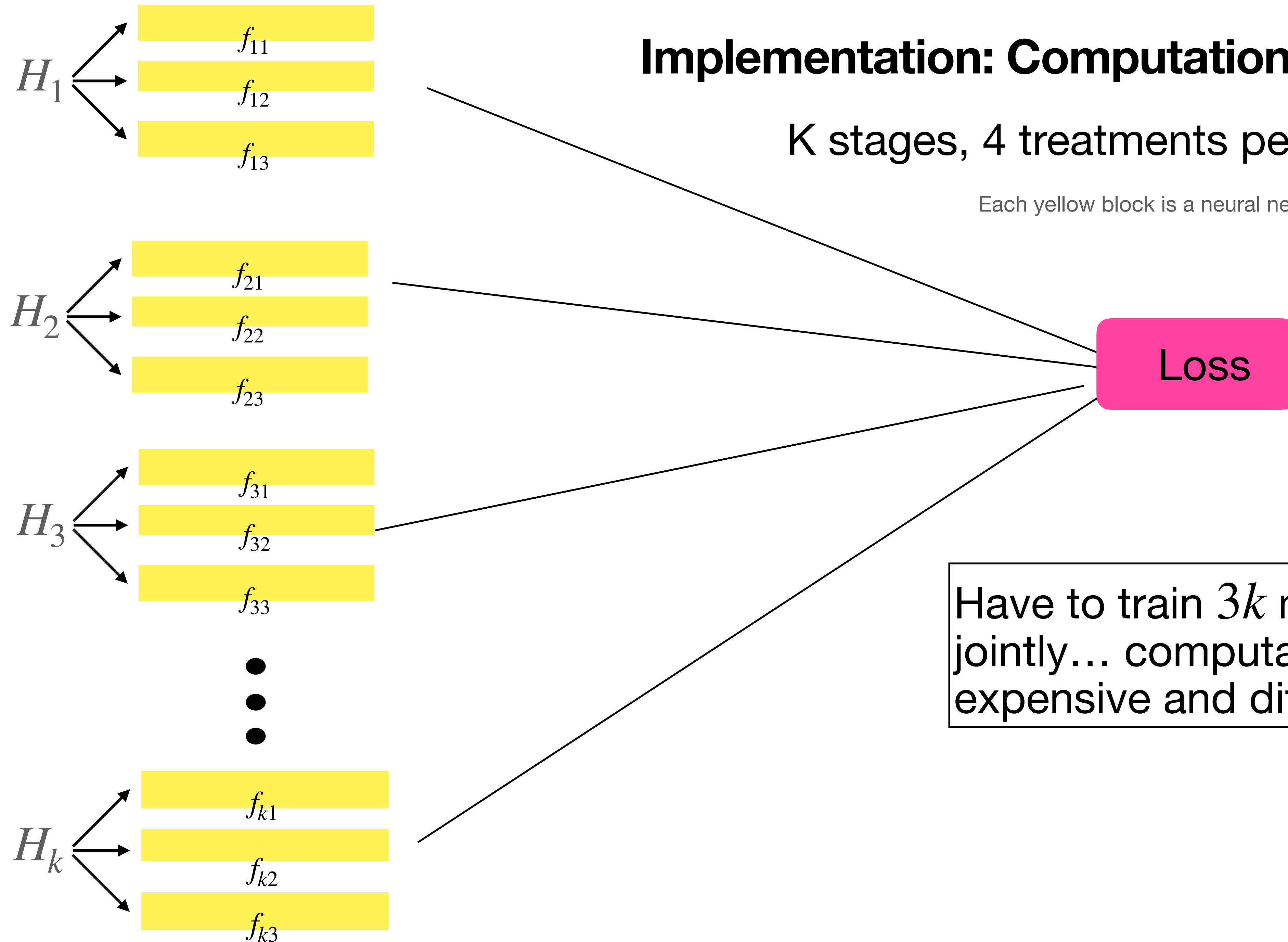
Each yellow block is a neural network



Implementation: Computational challenge

K stages, 4 treatments per stage

Each yellow block is a neural network

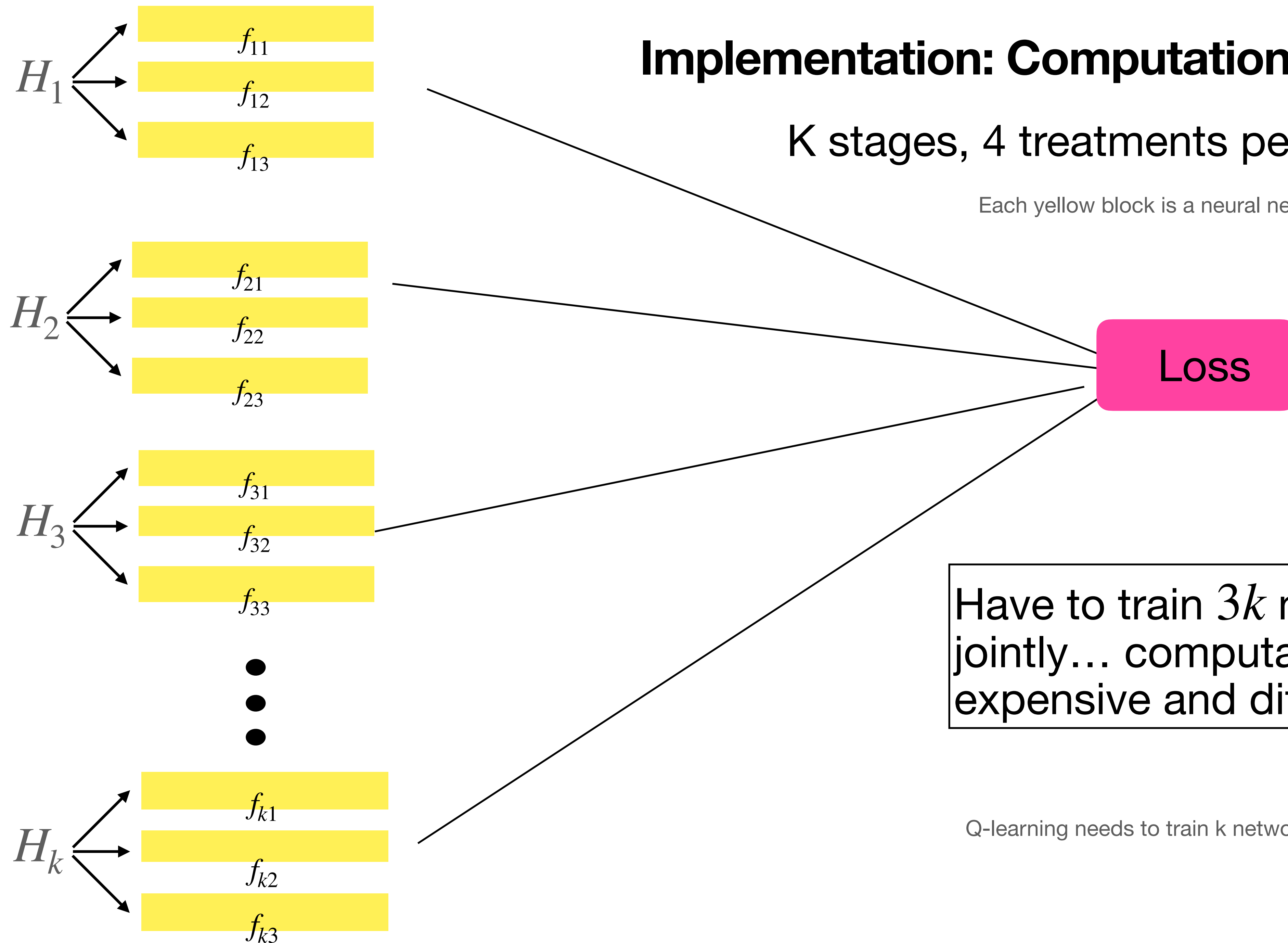


Have to train $3k$ networks jointly... computationally expensive and difficult to tune

Implementation: Computational challenge

K stages, 4 treatments per stage

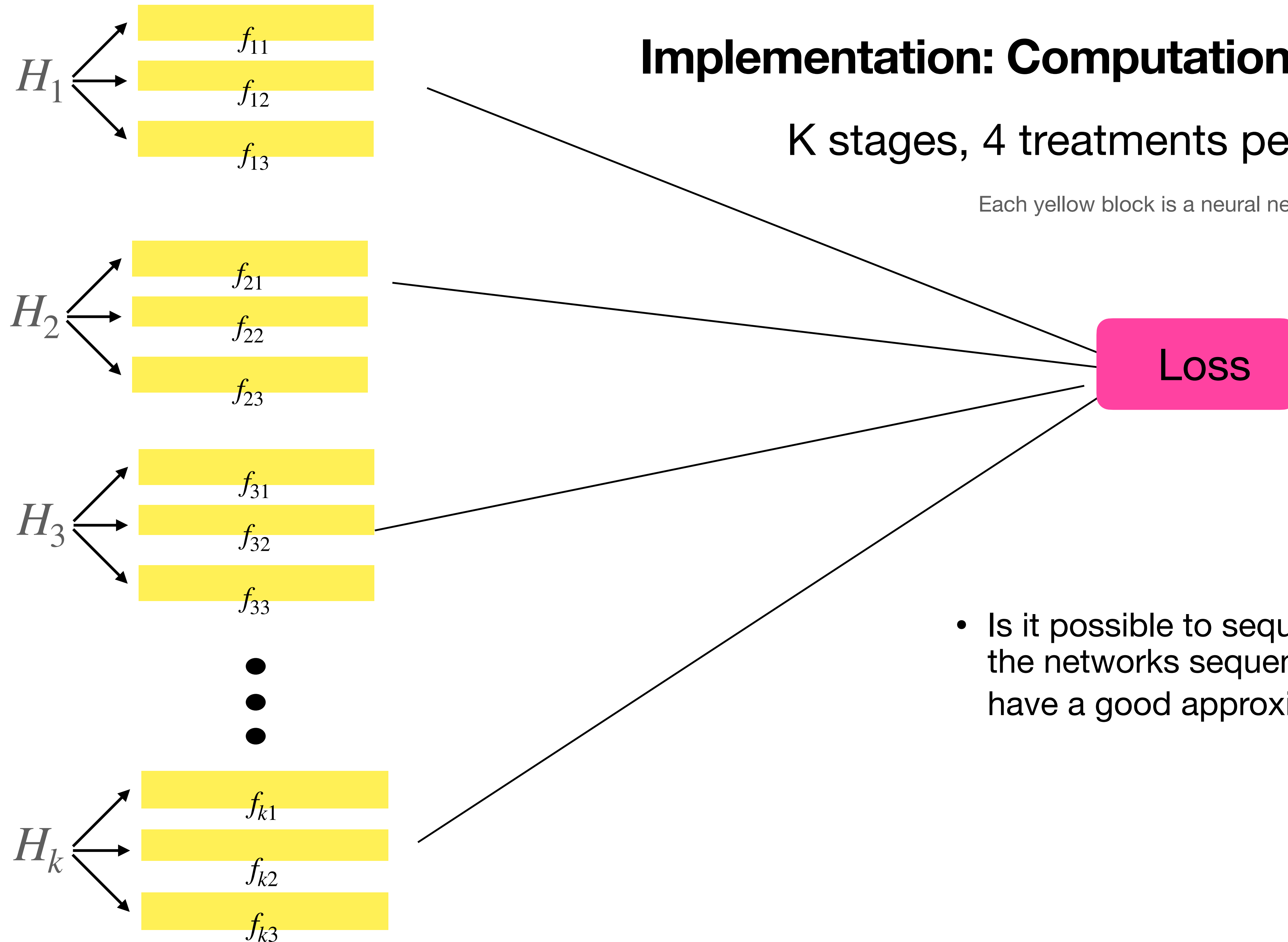
Each yellow block is a neural network



Implementation: Computational challenge

K stages, 4 treatments per stage

Each yellow block is a neural network

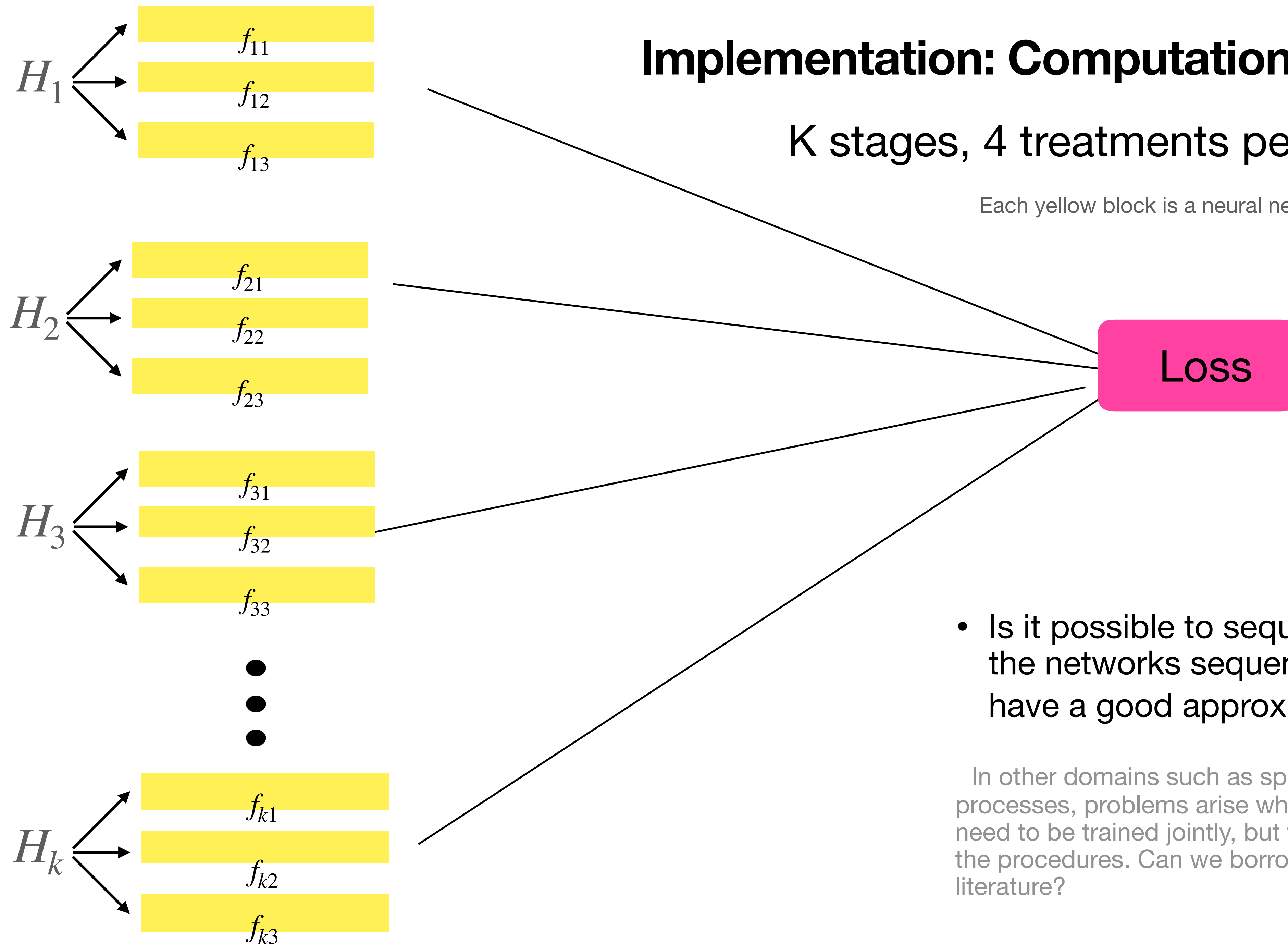


- Is it possible to sequentialize: train the networks sequentially but still have a good approximation of π^* ?

Implementation: Computational challenge

K stages, 4 treatments per stage

Each yellow block is a neural network



- Is it possible to sequentialize: train the networks sequentially but still have a good approximation of π^* ?

In other domains such as spatio-temporal processes, problems arise where neural networks need to be trained jointly, but they sequentialize the procedures. Can we borrow ideas from that literature?

Optimization

Optimization

Optimization

- Currently using stochastic gradient descent (SGD) for optimization — too general. We have a specific problem — can we tailor an optimization method?

Optimization

- Currently using stochastic gradient descent (SGD) for optimization — too general. We have a specific problem — can we tailor an optimization method?
- Feng et al. (2022) used similar loss function for another machine learning problem called ‘maximum score estimation’, and tailored an optimization method for their problem. Can we do something similar?

Optimization

- Currently using stochastic gradient descent (SGD) for optimization — too general. We have a specific problem — can we tailor an optimization method?
- Feng et al. (2022) used similar loss function for another machine learning problem called ‘maximum score estimation’, and tailored an optimization method for their problem. Can we do something similar?

Will require analysis of the optimization landscape

Optimization landscape

1. Nguyen et al., 2017 and 2019
2. Laha et al., 2022

Optimization landscape

1. Nguyen et al., 2017 and 2019
2. Laha et al., 2022

Optimization landscape

- Neural network classifiers:

Optimization landscape

- Neural network classifiers:
Existing deep learning results: can be used¹.

1. Nguyen et al., 2017 and 2019

2. Laha et al., 2022

Optimization landscape

- Neural network classifiers:

Existing deep learning results: can be used¹.

Challenges: loss non-standard, existing results not directly applicable

1. Nguyen et al., 2017 and 2019

2. Laha et al., 2022

Optimization landscape

K=1, 3 treatments, one covariate ($S_1 \in \mathbb{R}$), linear classifier

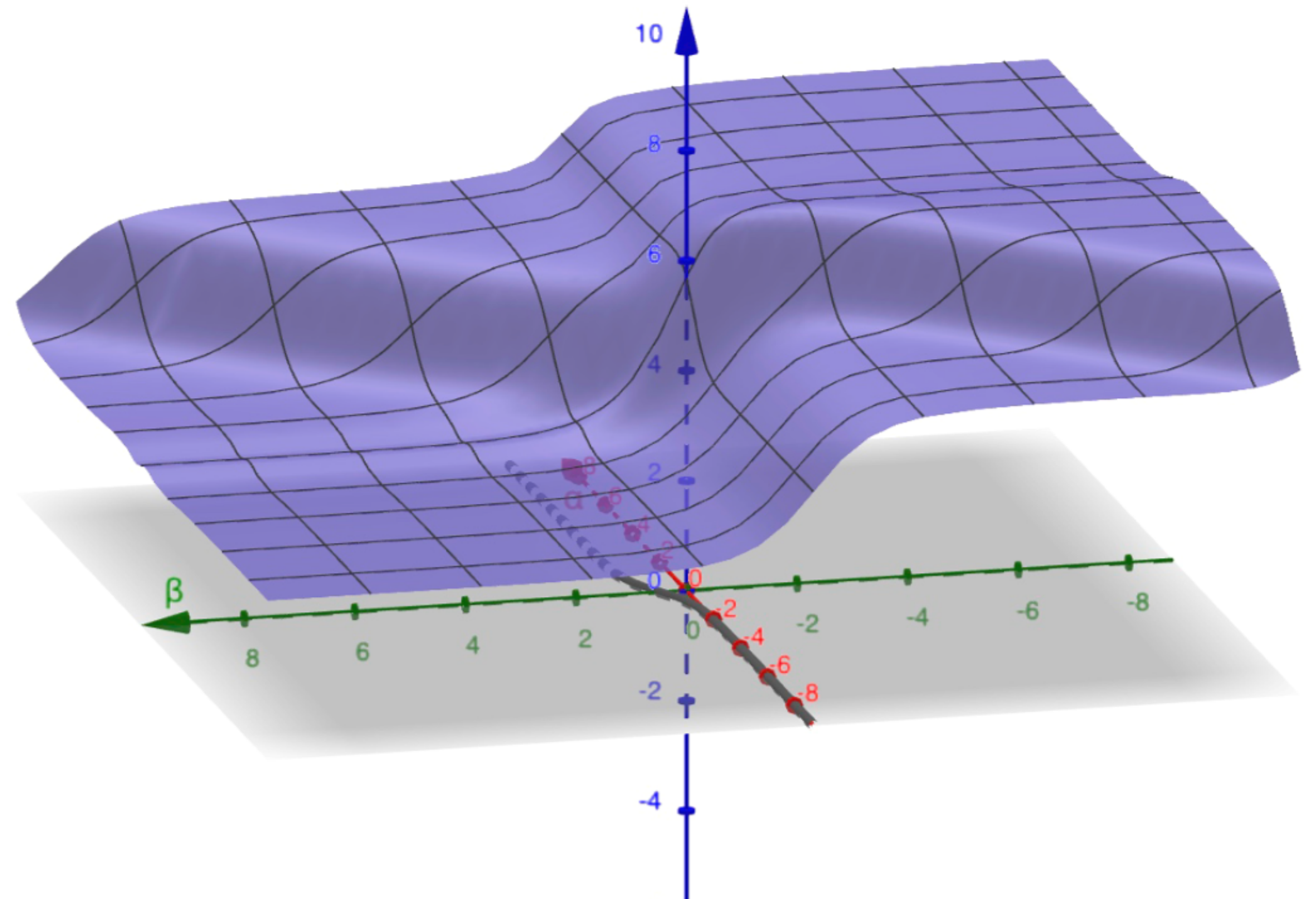
- Neural network classifiers:

Existing deep learning results: can be used¹.

Challenges: loss non-standard, existing results not directly applicable

- Linear classifiers:

optimization surface — specific properties: No local minima + regions with small gradient²



1. Nguyen et al., 2017 and 2019

2. Laha et al., 2022

Skills you will learn

Skills you will learn

Skills you will learn

DTR

Skills you will learn

DTR

PyTorch

Skills you will learn

DTR

PyTorch

Working with deep neural nets

Skills you will learn

DTR

PyTorch

Working with deep neural nets

Convergence of optimization-methods for non-convex problems¹

Skills you will learn

DTR

PyTorch

Working with deep neural nets

Convergence of optimization-methods for non-convex problems¹

Outline

- Example: sepsis
- Problem formulation
- Proposed method
- Open questions

A. Implementation and optimization

B. Regret decay rate

C. Doubly robust learning

Open questions

Regret decay

Open questions

Regret decay

- Regret: $V^{\pi^*} - V^{\hat{\pi}}$ measures how well we approximated π^* using $\hat{\pi}$.

Open questions

Regret decay

- Regret: $V^{\pi^*} - V^{\hat{\pi}}$ measures how well we approximated π^* using $\hat{\pi}$.

What is the rate of decay of regret?

Open questions

Regret decay

- Regret: $V^{\pi^*} - V^{\hat{\pi}}$ measures how well we approximated π^* using $\hat{\pi}$.

What is the rate of decay of regret?

Probably won't be very different from the 2-treatments case (Laha et al., 2022).

Open questions

Regret decay

- Regret: $V^{\pi^*} - V^{\hat{\pi}}$ measures how well we approximated π^* using $\hat{\pi}$.

What is the rate of decay of regret?

Probably won't be very different from the 2-treatments case (Laha et al., 2022).

Skills you will learn:

Open questions

Regret decay

- Regret: $V^{\pi^*} - V^{\hat{\pi}}$ measures how well we approximated π^* using $\hat{\pi}$.

What is the rate of decay of regret?

Probably won't be very different from the 2-treatments case (Laha et al., 2022).

Skills you will learn:

1. DTR

Open questions

Regret decay

- Regret: $V^{\pi^*} - V^{\hat{\pi}}$ measures how well we approximated π^* using $\hat{\pi}$.

What is the rate of decay of regret?

Probably won't be very different from the 2-treatments case (Laha et al., 2022).

Skills you will learn:

1. DTR
2. Empirical risk minimization theory

Open questions

Regret decay

- Regret: $V^{\pi^*} - V^{\hat{\pi}}$ measures how well we approximated π^* using $\hat{\pi}$.

What is the rate of decay of regret?

Probably won't be very different from the 2-treatments case (Laha et al., 2022).

Skills you will learn:

1. DTR
2. Empirical risk minimization theory
3. Some theory on multcategory classification

Open questions

Regret decay

- Regret: $V^{\pi^*} - V^{\hat{\pi}}$ measures how well we approximated π^* using $\hat{\pi}$.

What is the rate of decay of regret?

Probably won't be very different from the 2-treatments case (Laha et al., 2022).

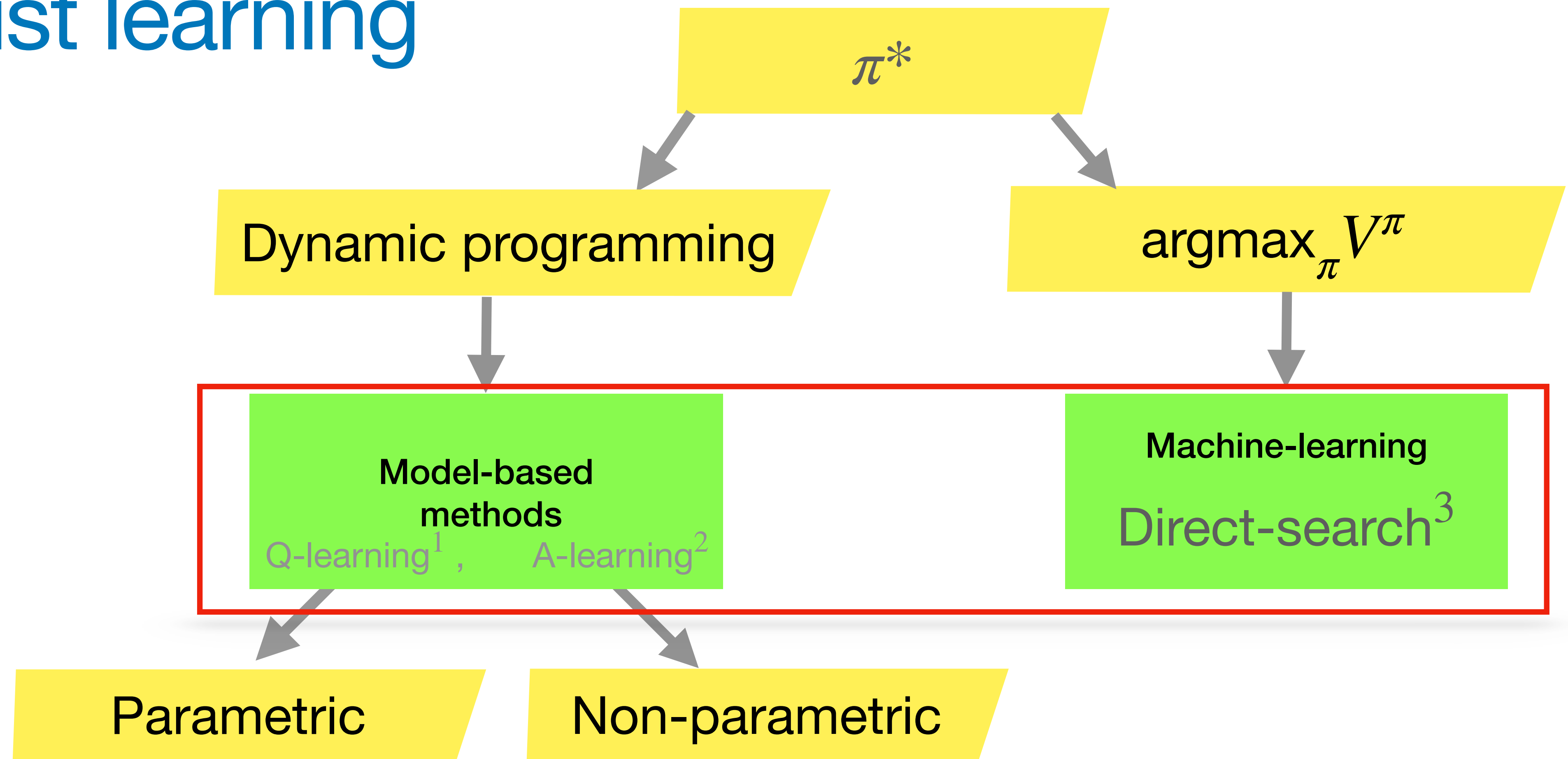
Skills you will learn:

1. DTR
2. Empirical risk minimization theory
3. Some theory on multcategory classification
4. Some theory on policy learning in offline RL

Outline

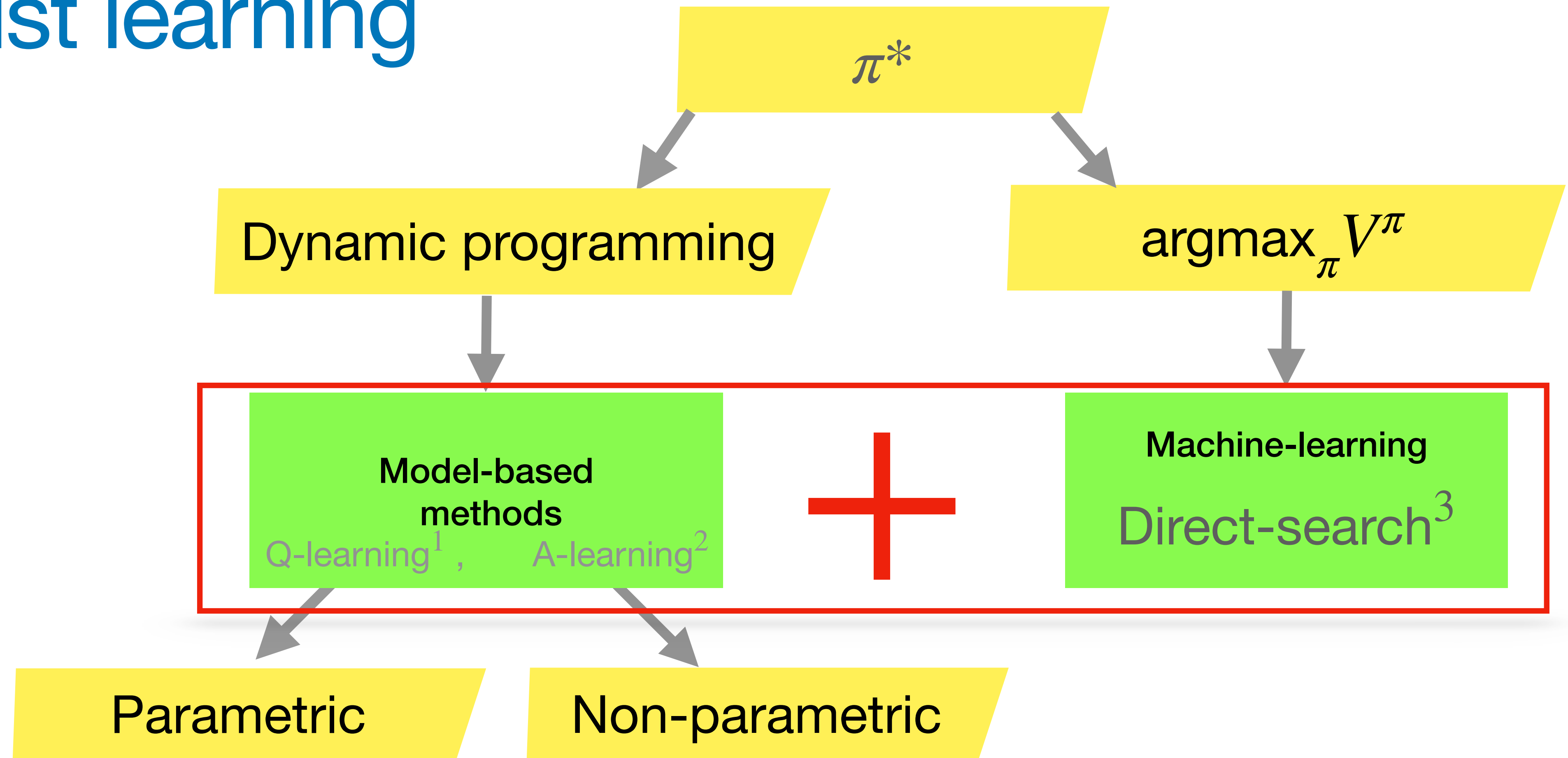
- Example: sepsis
- Problem formulation
- Proposed method
- Open questions
 - A. Implementation and optimization
 - B. Regret decay rate
 - C. Doubly robust learning

Doubly robust learning



1. Watkins, 1989; Schulte et al. 2014
2. Murphy, 2003; Robins, 2004
3. Zhao et al. 2012; 2015, Laha et al. 2023

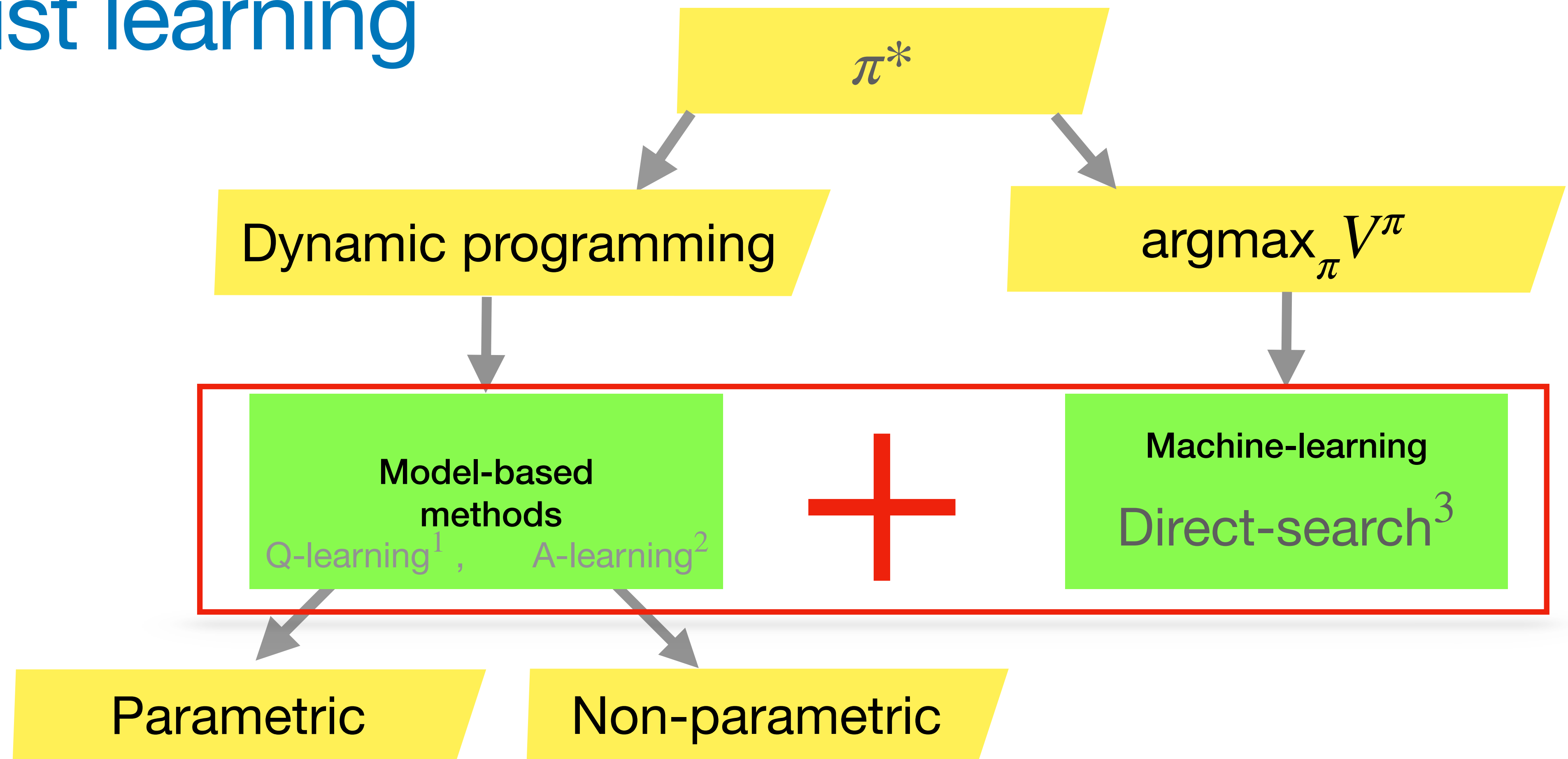
Doubly robust learning



Hybrid method (idea taken from offline RL)

1. Watkins, 1989; Schulte et al. 2014
2. Murphy, 2003; Robins, 2004
3. Zhao et al. 2012; 2015, Laha et al. 2023

Doubly robust learning

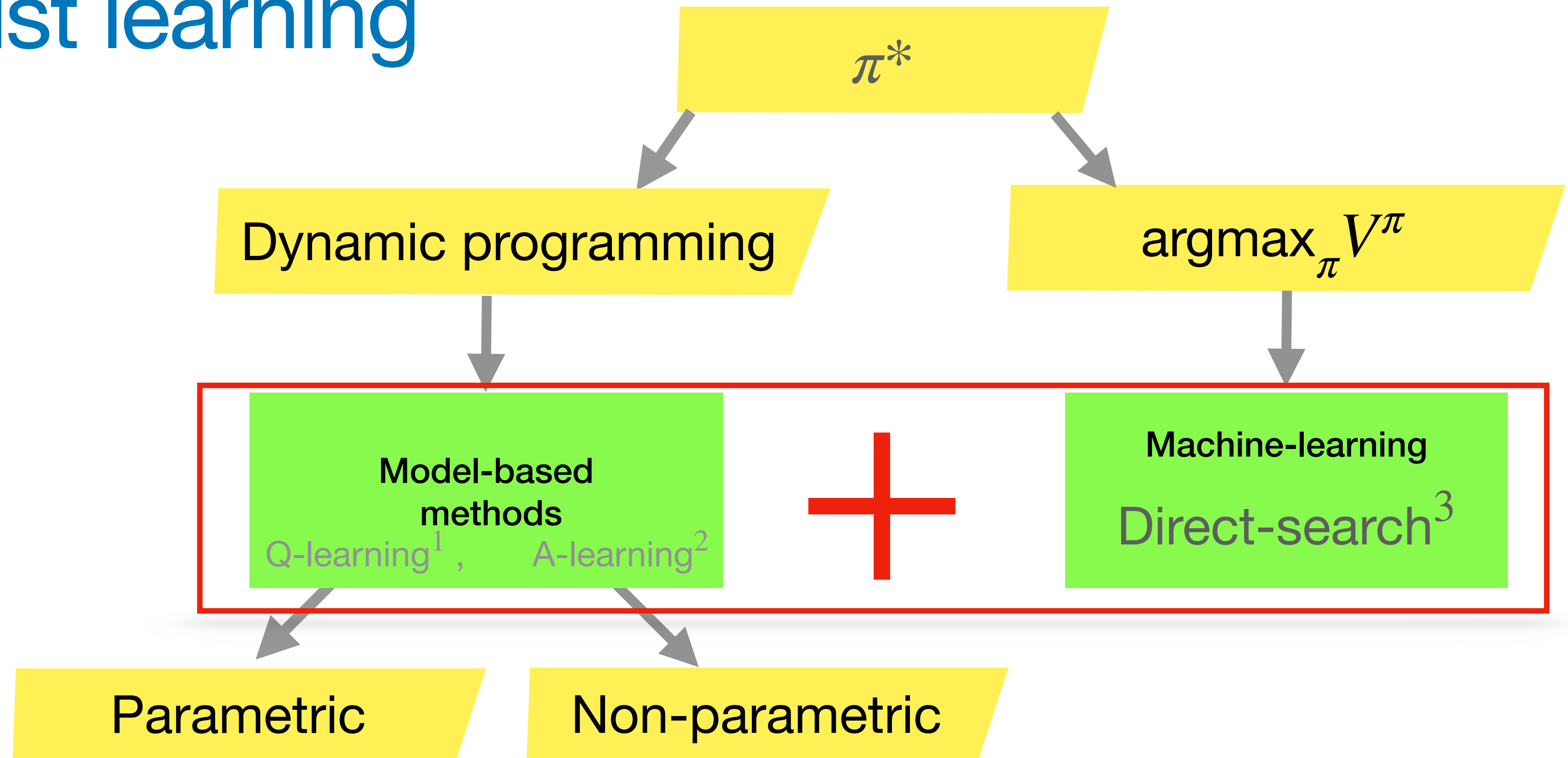


Hybrid method (idea taken from offline RL)

If either the Q-learning model assumptions or the estimation of treatment assignment probabilities correct, then π^* consistently estimated

1. Watkins, 1989; Schulte et al. 2014
2. Murphy, 2003; Robins, 2004
3. Zhao et al. 2012; 2015, Laha et al. 2023

Doubly robust learning



Hybrid method (idea taken from offline RL)

If either the Q-learning model assumptions or the estimation of treatment assignment probabilities correct, then π^* consistently estimated

1. Watkins, 1989; Schulte et al. 2014
2. Murphy, 2003; Robins, 2004
3. Zhao et al. 2012; 2015, Laha et al. 2023

Open questions

Open questions

1. I already have the method, but same questions on implementation

Open questions

1. I already have the method, but same questions on implementation
2. regret decay: \sqrt{n} – consistent

Open questions

1. I already have the method, but same questions on implementation
2. regret decay: \sqrt{n} – consistent

Skills you will learn:

Open questions

1. I already have the method, but same questions on implementation
2. regret decay: \sqrt{n} —consistent

Skills you will learn:

1. DTR

Open questions

1. I already have the method, but same questions on implementation
2. regret decay: \sqrt{n} – consistent

Skills you will learn:

1. DTR
2. Q-learning

Open questions

1. I already have the method, but same questions on implementation
2. regret decay: \sqrt{n} – consistent

Skills you will learn:

1. DTR
2. Q-learning
3. doubly robust offline RL

Open questions

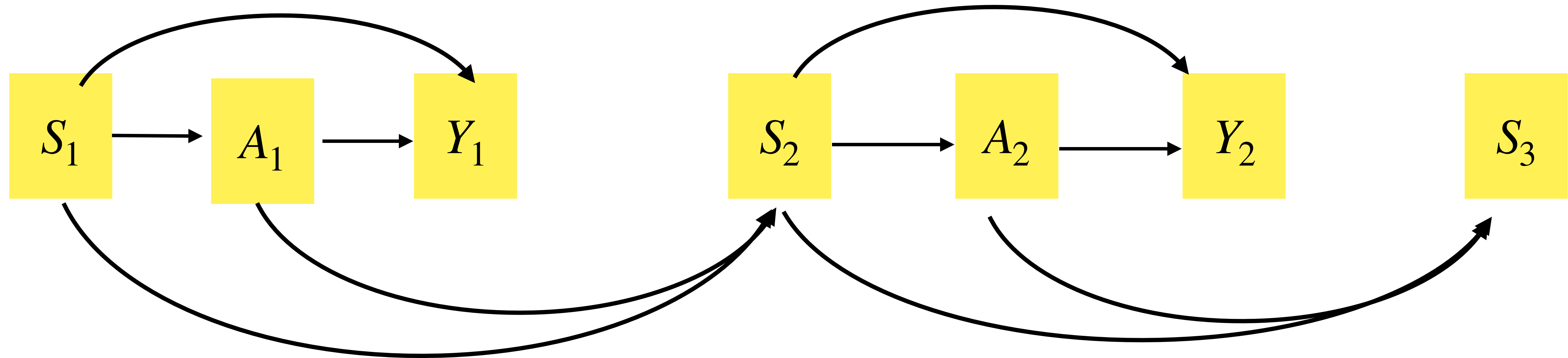
1. I already have the method, but same questions on implementation
2. regret decay: \sqrt{n} – consistent

Skills you will learn:

1. DTR
2. Q-learning
3. doubly robust offline RL
4. Some doubly robust literature in causal inference

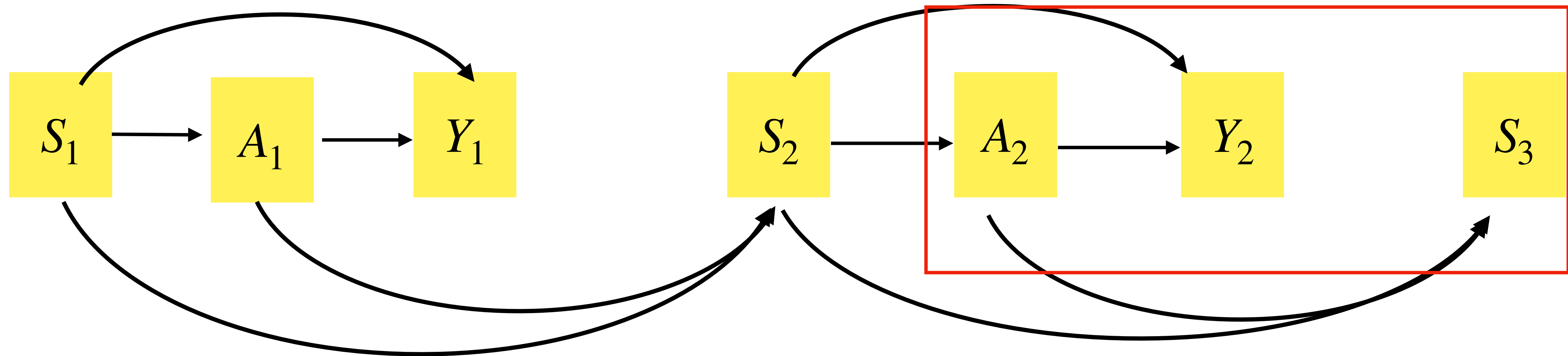
Full reinforcement learning

We do not make Markov decision process (MDP) assumption



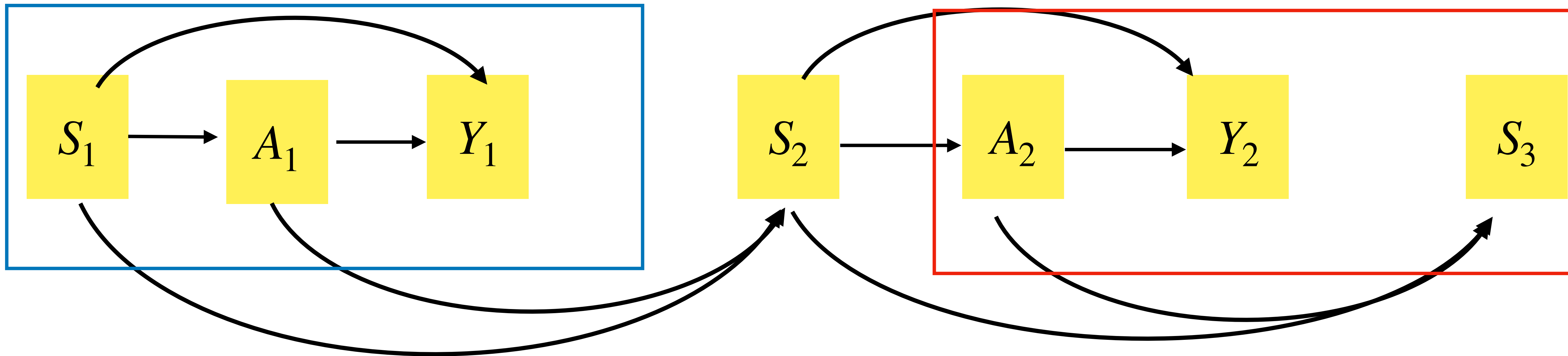
Full reinforcement learning

We do not make Markov decision process (MDP) assumption

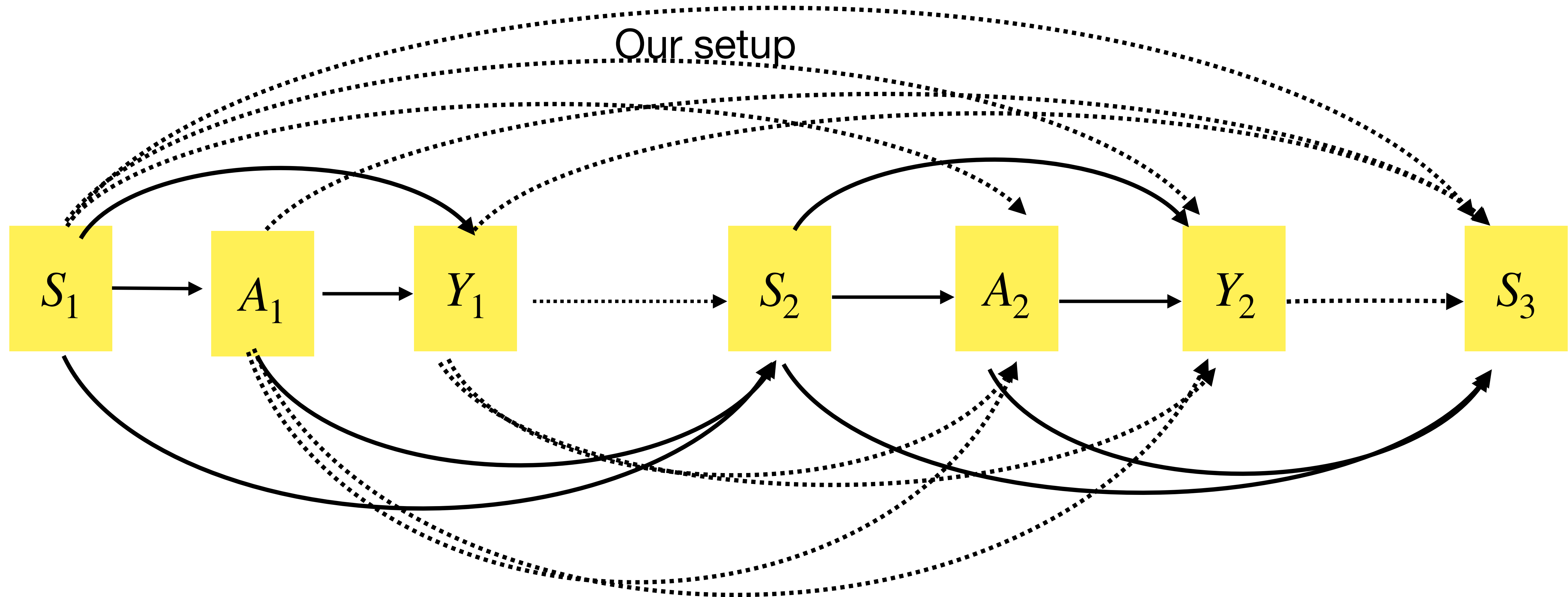


Full reinforcement learning

We do not make Markov decision process (MDP) assumption

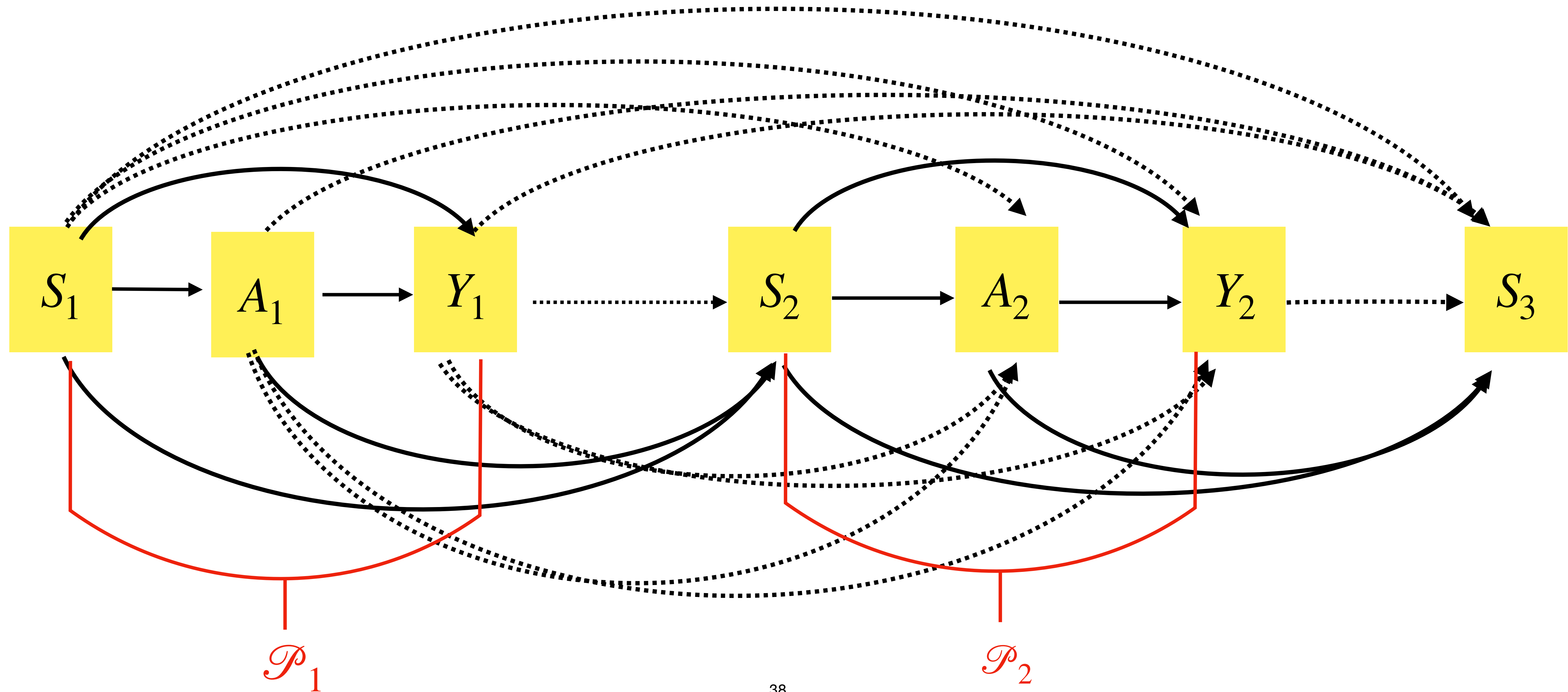


Full reinforcement learning

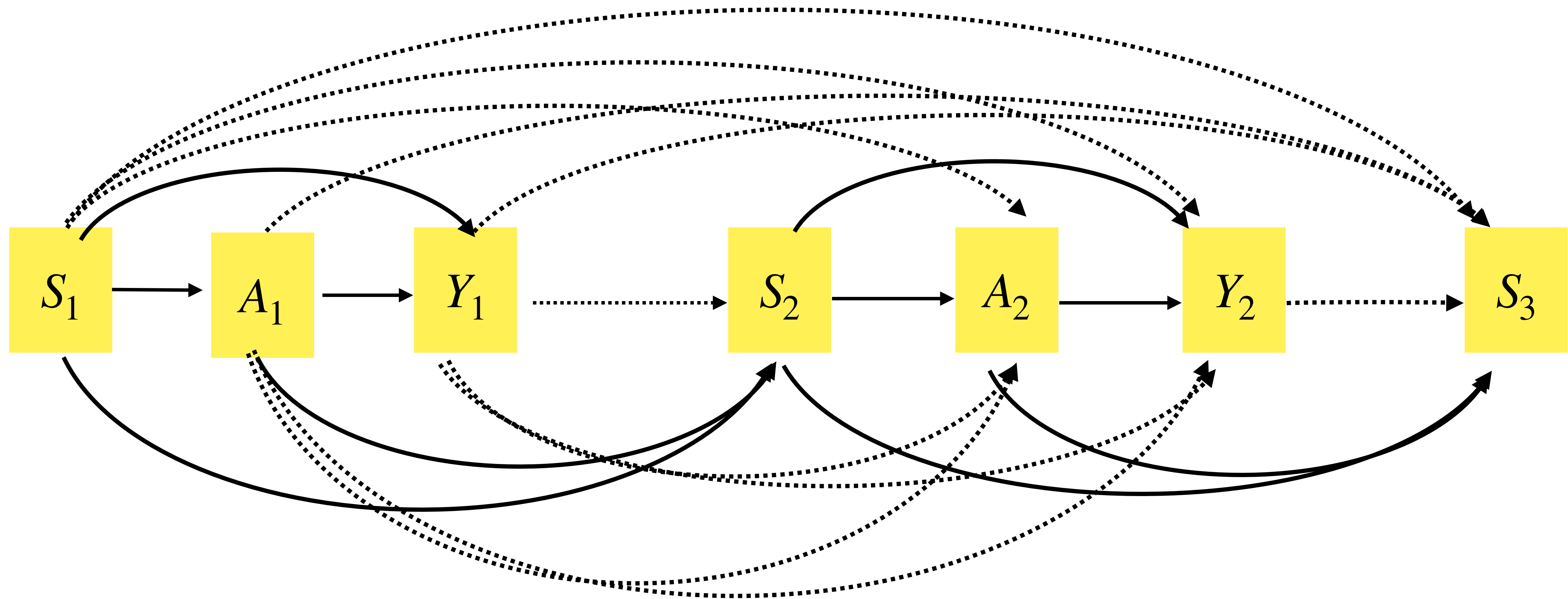


Full reinforcement learning

No stationarity



Full reinforcement learning



Set-up

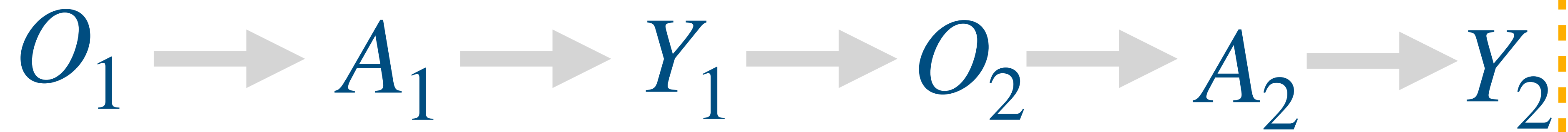


Set-up



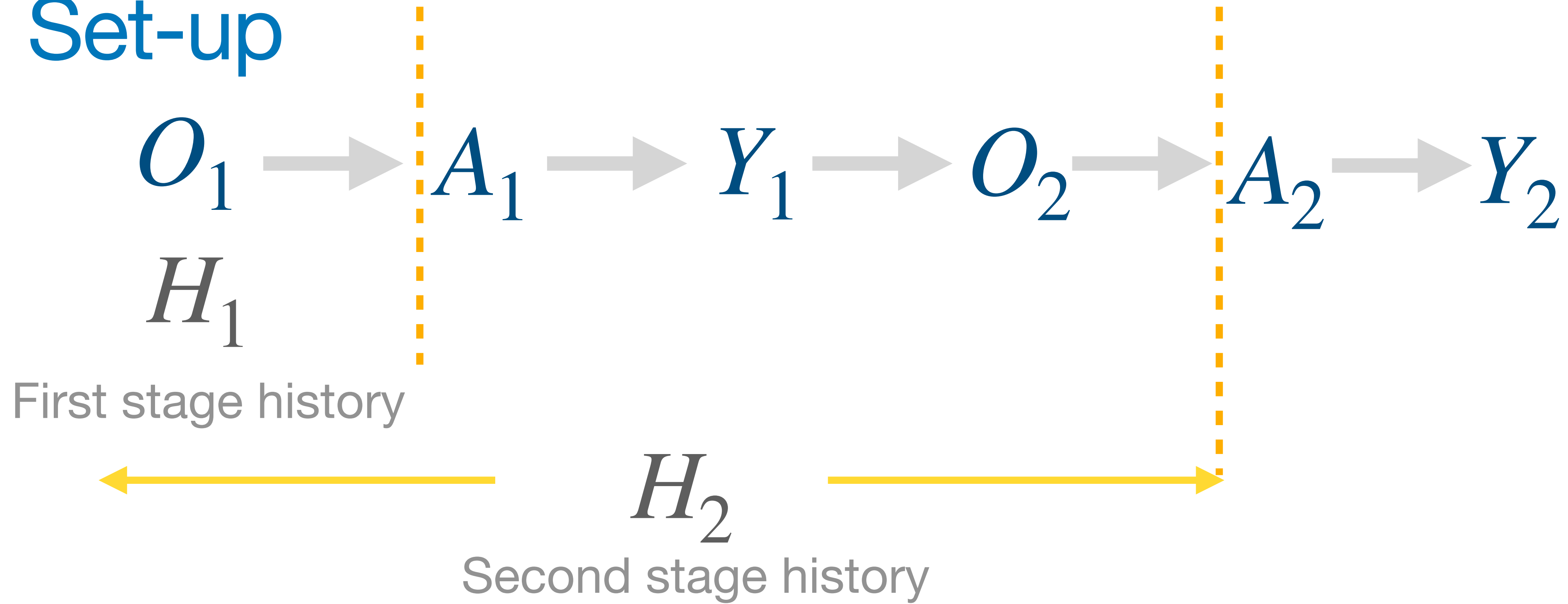
$K=2$

Set-up

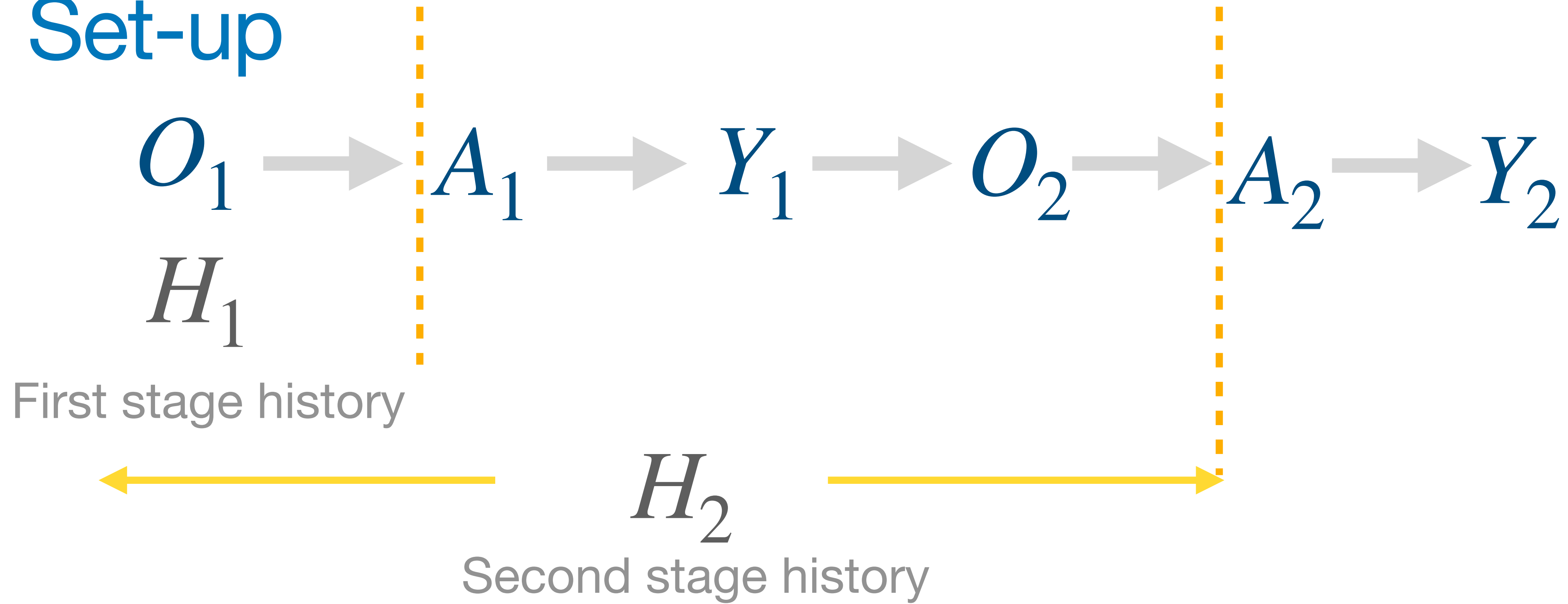


$K=2$

Set-up



Set-up



Treatment policy

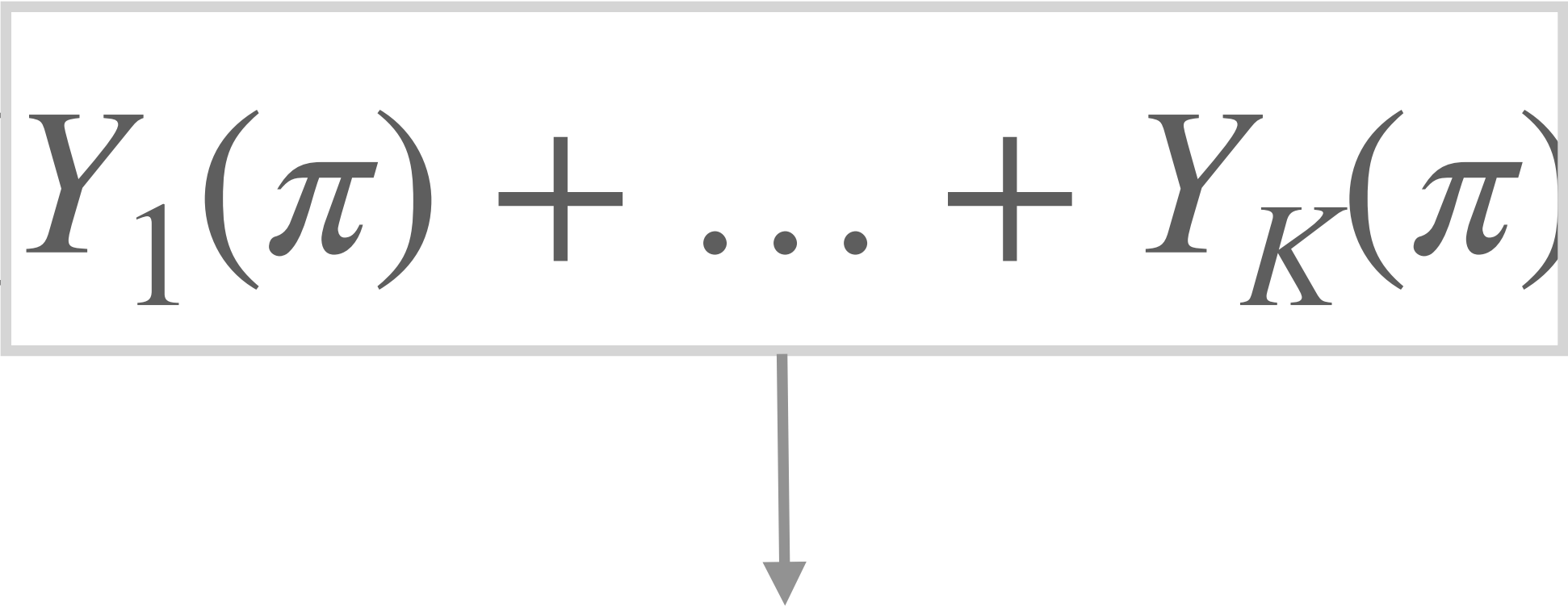
$$\pi = (\pi_1, \pi_2)$$

Value function estimation

$$V^\pi = \mathbb{E}[Y_1(\pi) + \dots + Y_K(\pi)]$$

Optimal treatment assignment $\pi^* = \operatorname{argmax}_\pi V^\pi$

Value function estimation


$$V^\pi = \mathbb{E}[Y_1(\pi) + \dots + Y_K(\pi)]$$


Potential outcomes

Optimal treatment assignment $\pi^* = \operatorname{argmax}_\pi V^\pi$

Value function estimation

Under standard identifiability assumptions*,

$$V^\pi = \mathbb{E} \left[(Y_1 + \dots + Y_K) \frac{\pi_1(A_1 | H_1) \dots \pi_K(A_K | H_K)}{\pi_{b,1}(A_1 | H_1) \dots \pi_{b,K}(A_K | H_K)} \right]$$


$\pi_{b,k}$'s behavior policy: ratio called inverse probability weights

Optimal treatment assignment $\pi^* = \operatorname{argmax}_\pi V^\pi$

*Orellana et al., 2010

Value function estimation

Under standard identifiability assumptions*,

$$V^\pi \approx \mathbb{P}_n \left[(Y_1 + \dots + Y_K) \frac{\pi_1(A_1 | H_1) \dots \pi_K(A_K | H_K)}{\pi_{b,1}(A_1 | H_1) \dots \pi_{b,K}(A_K | H_K)} \right]$$

\mathbb{P}_n : empirical distribution function

Optimal treatment assignment $\pi^* = \operatorname{argmax}_\pi V^\pi$

*Orellana et al., 2010

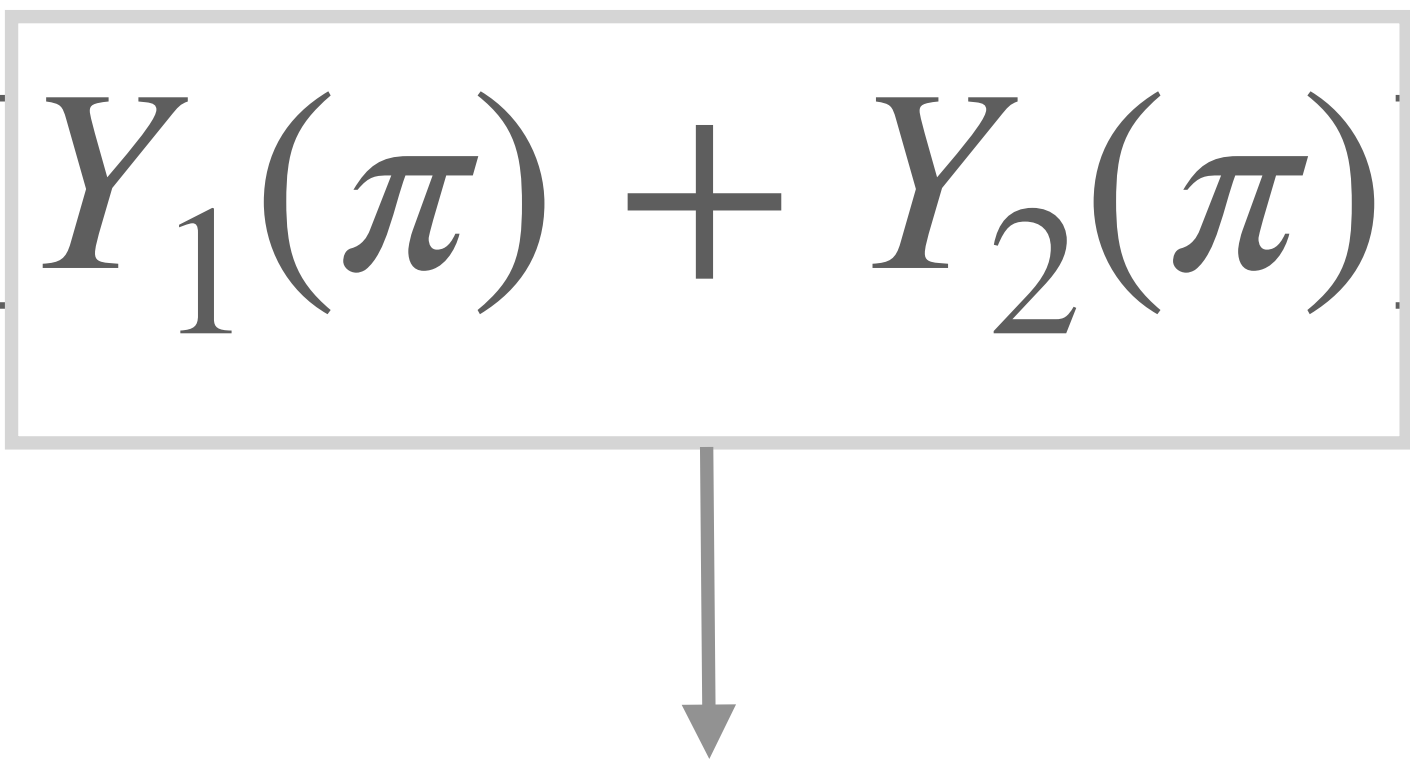
Value function

Optimal treatment policy $\pi^* = \operatorname{argmax}_{\pi} V^{\pi}$

$$V^{\pi} = \mathbb{E}[Y_1(\pi) + Y_2(\pi)]$$

Value function

Optimal treatment policy $\pi^* = \operatorname{argmax}_{\pi} V^{\pi}$

$$V^{\pi} = \mathbb{E}[Y_1(\pi) + Y_2(\pi)]$$


Potential outcomes

Value function

Optimal treatment policy $\pi^* = \operatorname{argmax}_{\pi} V^{\pi}$

Under standard identifiability assumptions*,

$$V^{\pi} = \mathbb{E} \left[(Y_1 + Y_2) \frac{1\{\pi_1(H_1) = A_1\} 1\{\pi_2(H_2) = A_2\}}{P(A_1 | H_1) P(A_2 | H_2)} \right]$$

↓
observed random variables

Value function

Optimal treatment policy $\pi^* = \operatorname{argmax}_{\pi} V^{\pi}$

Under standard identifiability assumptions*,

$$V^{\pi} \approx \frac{1}{n} \sum_{i=1}^n \left((Y_{1i} + Y_{2i}) \frac{1\{\pi_1(H_{1i}) = A_{1i}\} 1\{\pi_2(H_{2i}) = A_{2i}\}}{P(A_{1i} | H_{1i}) P(A_{2i} | H_{2i})} \right)$$

*Orellana et al., 2010

Value function

Optimal treatment policy $\pi^* = \operatorname{argmax}_{\pi} V^{\pi}$

Under standard identifiability assumptions*,

$$V^{\pi} \approx \frac{1}{n} \sum_{i=1}^n \left((Y_{1i} + Y_{2i}) \frac{1\{\pi_1(H_{1i}) = A_{1i}\} 1\{\pi_2(H_{2i}) = A_{2i}\}}{P(A_{1i} | H_{1i}) P(A_{2i} | H_{2i})} \right)$$

Maximize V^{π} over a. Class of policies

*Orellana et al., 2010

Value function

Optimal treatment policy $\pi^* = \operatorname{argmax}_{\pi} V^{\pi}$

Under standard identifiability assumptions*,

$$V^{\pi} \approx \frac{1}{n} \sum_{i=1}^n \left((Y_{1i} + Y_{2i}) \frac{1\{\pi_1(H_{1i}) = A_{1i}\} 1\{\pi_2(H_{2i}) = A_{2i}\}}{P(A_{1i} | H_{1i}) P(A_{2i} | H_{2i})} \right)$$

Discontinuous + non, convex


Direct optimization not computationally feasible

*Orellana et al., 2010

Shortcomings of previous method

- $\min_{f:H \mapsto \mathbb{R}^4} E \left[C(H_1, Y_1) \times 1[\operatorname{argmax}(f(H_1)) \neq A_1] \right]$

Shortcomings of previous method

- $\min_{f:H \mapsto \mathbb{R}^4} E \left[C(H_1, Y_1) \times 1[\operatorname{argmax}(f(H_1)) \neq A_1] \right]$

$$C(H_1, Y_1) = \frac{Y_1}{P(A_1 | H_1)}$$

Shortcomings of previous method

- $\min_{f:H \mapsto \mathbb{R}^4} E \left[C(H_1, Y_1) \times 1[\operatorname{argmax}(f(H_1)) \neq A_1] \right]$



$$C(H_1, Y_1) = \frac{Y_1}{P(A_1 | H_1)}$$

If I don't know what doctors were thinking, need to model the probabilities

Shortcomings of previous method

- $\min_{f:H \mapsto \mathbb{R}^4} E \left[C(H_1, Y_1) \times 1[\operatorname{argmax}(f(H_1)) \neq A_1] \right]$



$$C(H_1, Y_1) = \frac{Y_1}{P(A_1 | H_1)}$$

If I don't know what doctors were thinking, need to model the probabilities


Bad estimation



$\hat{\pi}$ bad estimator of π^*

Shortcomings of previous method

- $\min_{f:H \mapsto \mathbb{R}^4} E \left[C(H_1, Y_1) \times 1[\operatorname{argmax}(f(H_1)) \neq A_1] \right]$


$$C(H_1, Y_1) = \frac{Y_1}{P(A_1 | H_1)}$$

If I don't know what doctors were thinking, need to model the probabilities

$P(A_1 | H_1)$ is small \implies the estimator of $C(H_1, A_1)$ can be highly variable

Loss function when stage $K = 1$

Classifiers for stage 1

Classifiers for stage 1



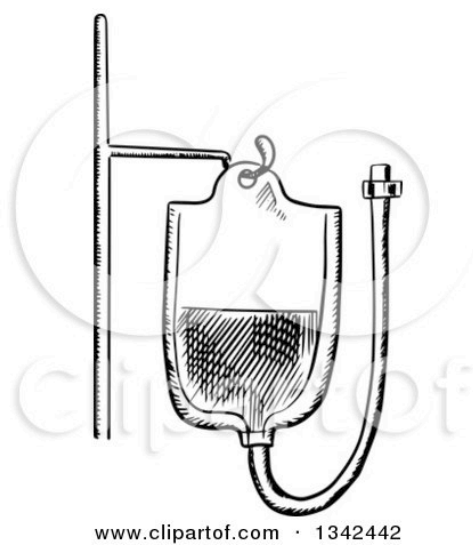
H_1

Classifiers for stage 1

Possible categories



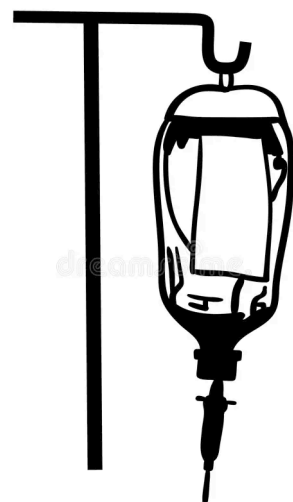
H_1



Low



Medium



High



No IV

Classifiers for stage 1

Possible categories

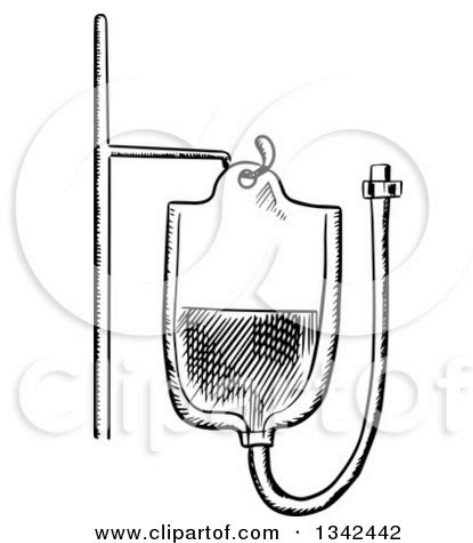


H_1

Classifier:

$$f = (f_1, \dots, f_4)$$

$$f_i : H_1 \mapsto \mathbb{R} \quad i = 1, \dots, 4$$



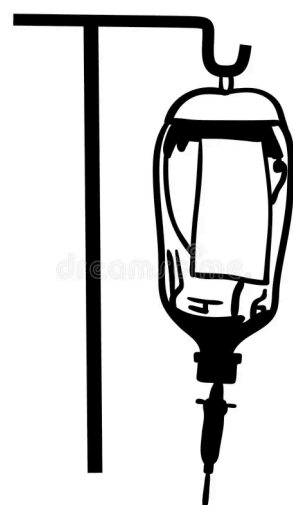
Low

$$f_1(H_1)$$



Medium

$$f_2(H_1)$$



High

$$f_3(H_1)$$



No IV

$$f_4(H_1)$$

Classifiers for stage 1

Possible categories

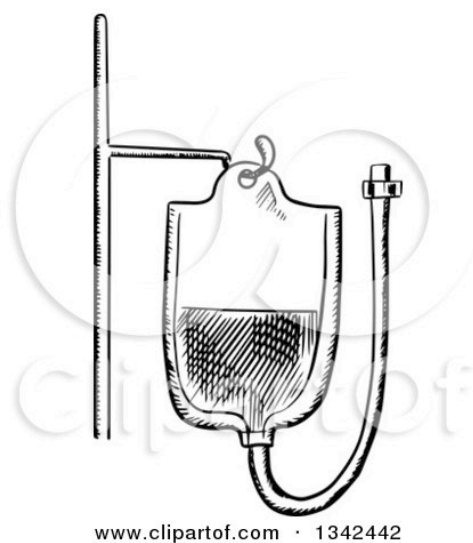


H_1

Classifier:

$$f = (f_1, \dots, f_4)$$

$$f_i : H_1 \mapsto \mathbb{R} \quad i = 1, \dots, 4$$



Low

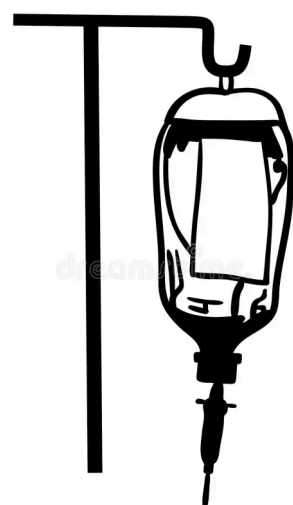
$$f_1(H_1)$$



Medium

$$f_2(H_1)$$

Maximum



High

$$f_3(H_1)$$



No IV

$$f_4(H_1)$$

Classifiers for stage 1

Possible categories

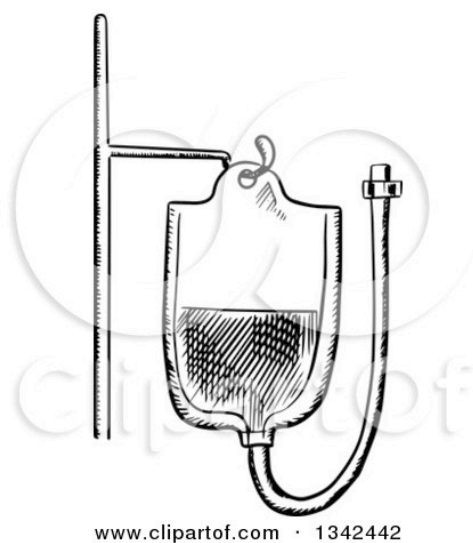


H_1

Classifier:

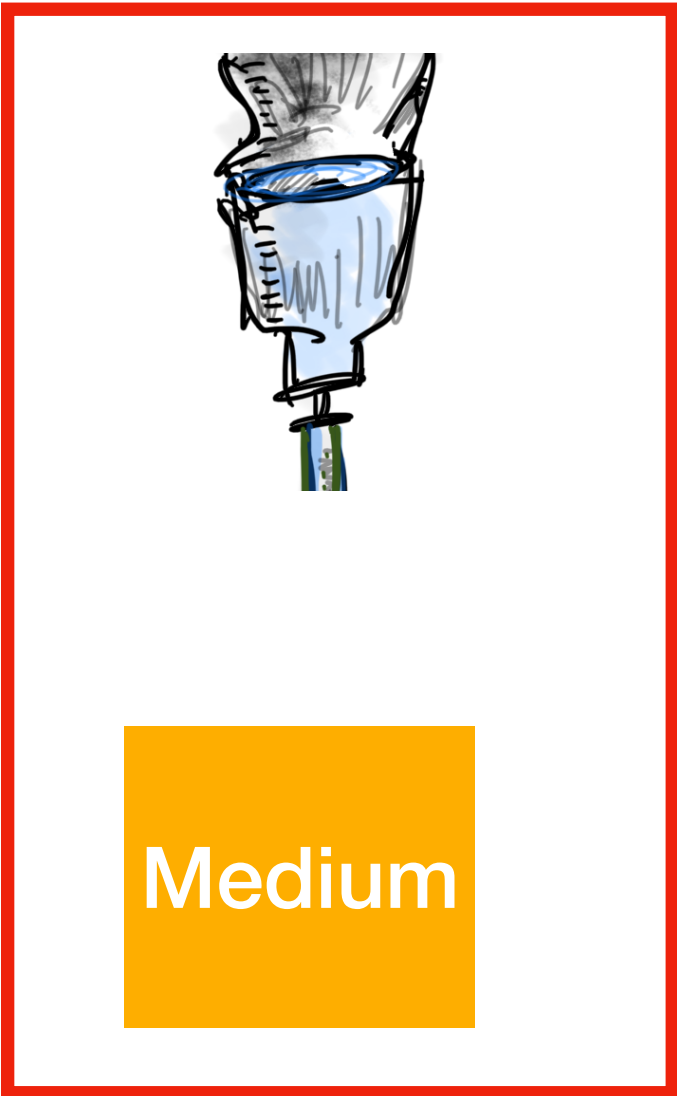
$$f = (f_1, \dots, f_4)$$

$$f_i : H_1 \mapsto \mathbb{R} \quad i = 1, \dots, 4$$



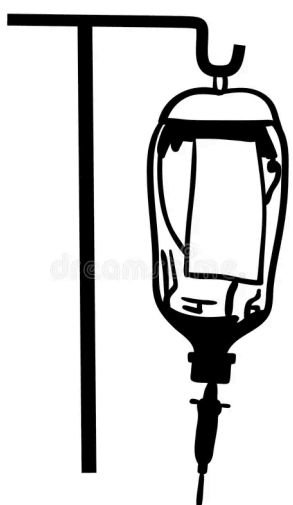
Low

$$f_1(H_1)$$



Medium

$$f_2(H_1)$$



High

$$f_3(H_1)$$



No IV

$$f_4(H_1)$$

Maximum

$$\pi_1(H_1) = \operatorname{argmax}_i f_i(H_1)$$

The loss function

Case $T = 1$

- $\max_{f: H_1 \mapsto \mathbb{R}^4} E \left[C(H_1, Y_1) \times 1[\operatorname{argmax}_i f_i(H_1) \neq A_1] \right]$

The loss function

Case $T = 1$

- $\max_{f:H_1 \mapsto \mathbb{R}^4} E [C(H_1, Y_1) \times 1[\operatorname{argmax}_i f_i(H_1) \neq A_1]]$



In practice search
over a smaller
class, currently we
consider neural
network classes

The loss function

Case $T = 1$

- $\max_{f: H_1 \mapsto \mathbb{R}^4} E \left[C(H_1, Y_1) \times 1[\operatorname{argmax}_i f_i(H_1) \neq A_1] \right]$



Depends on data

The loss function

Case $T = 1$

- $\max_{f:H_1 \mapsto \mathbb{R}^4} E \left[C(H_1, Y_1) \times 1[\operatorname{argmax}_i f_i(H_1) \neq A_1] \right]$

↓
Discontinuity

The loss function

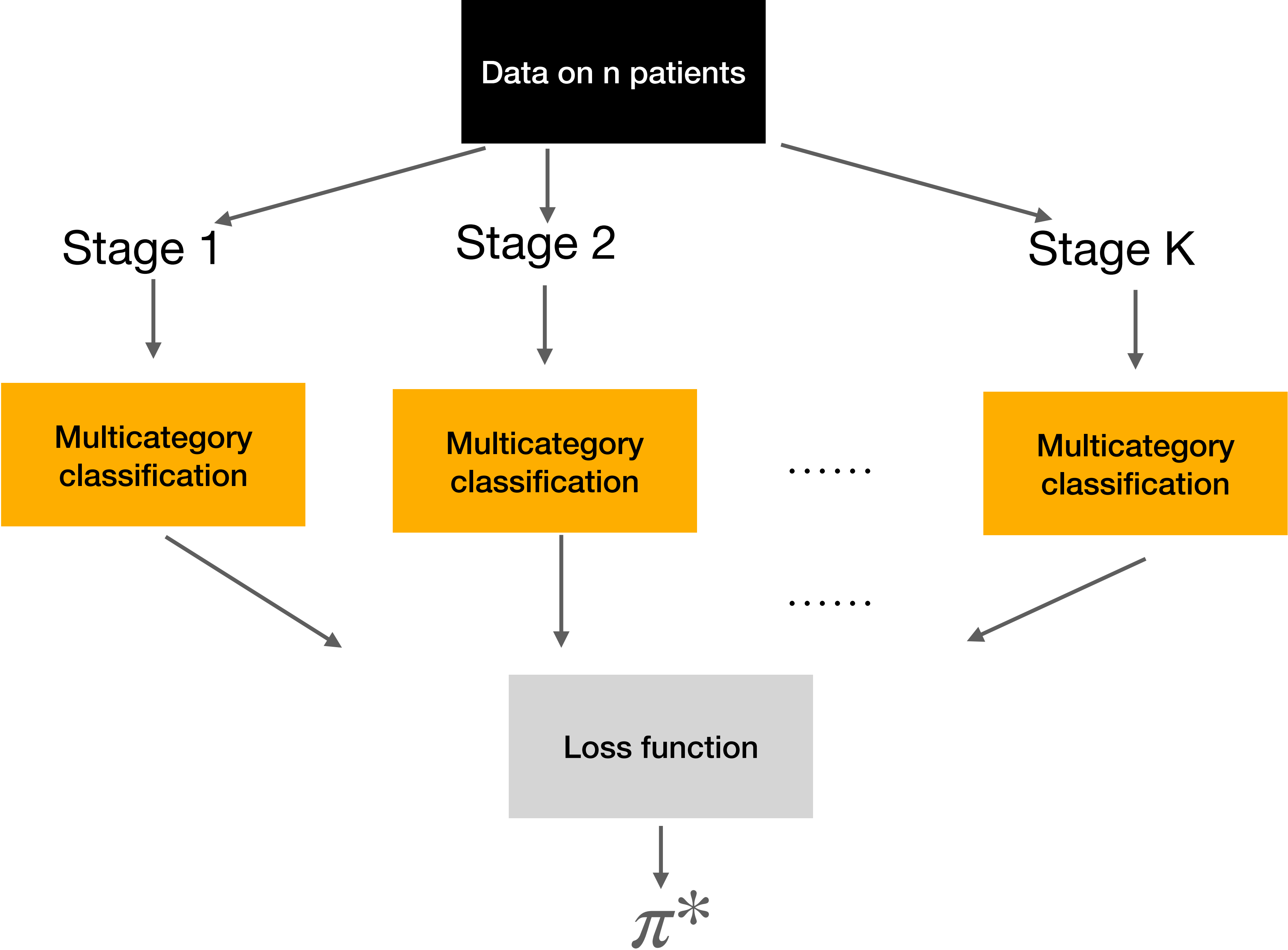
Case $T = 1$

- $\max_{f: H_1 \mapsto \mathbb{R}^4} E \left[C(H_1, Y_1) \times 1[\operatorname{argmax}_i f_i(H_1) \neq A_1] \right]$

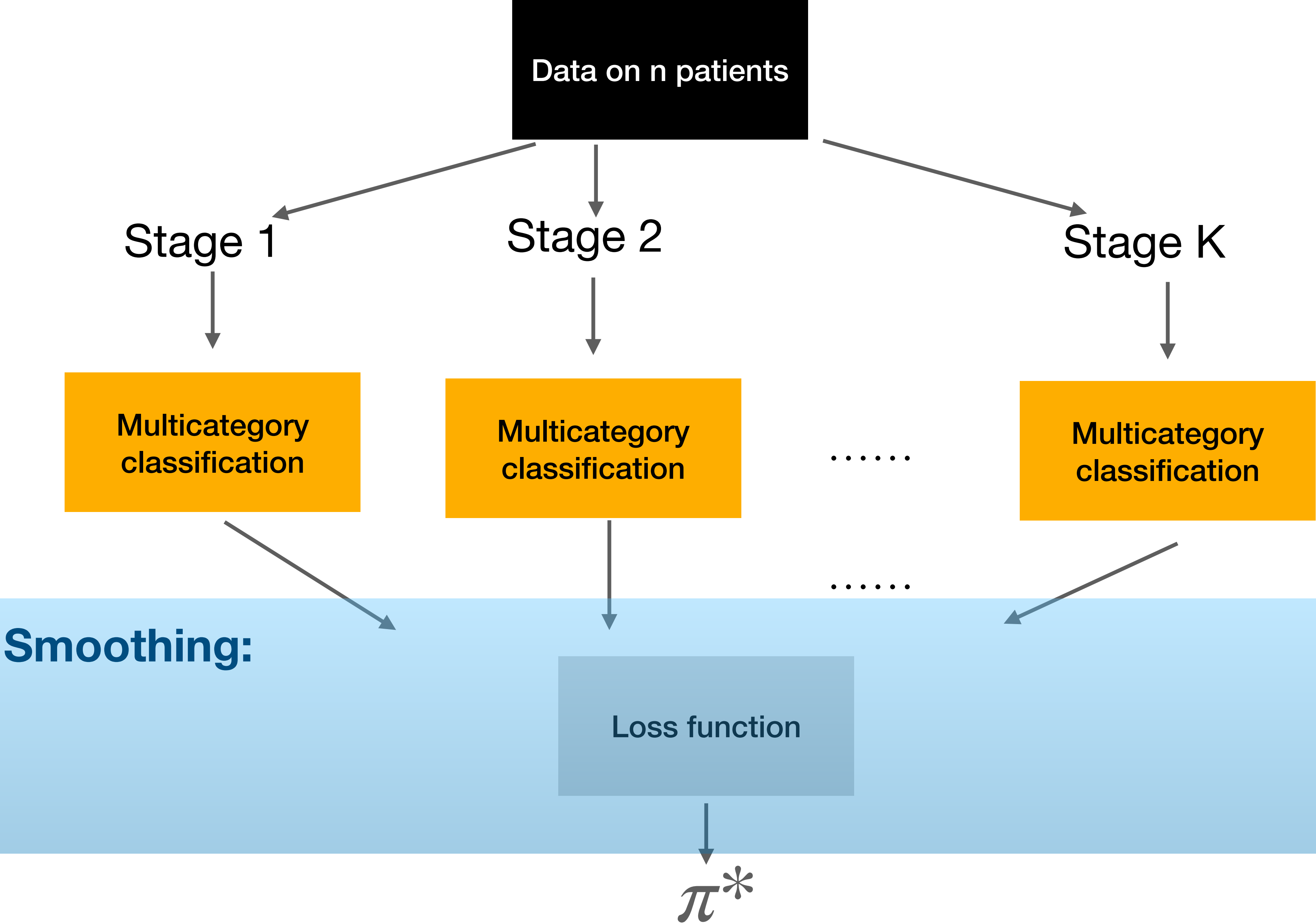
↓
Discontinuity

Our proposal: smooth out the sources for discontinuity at each step

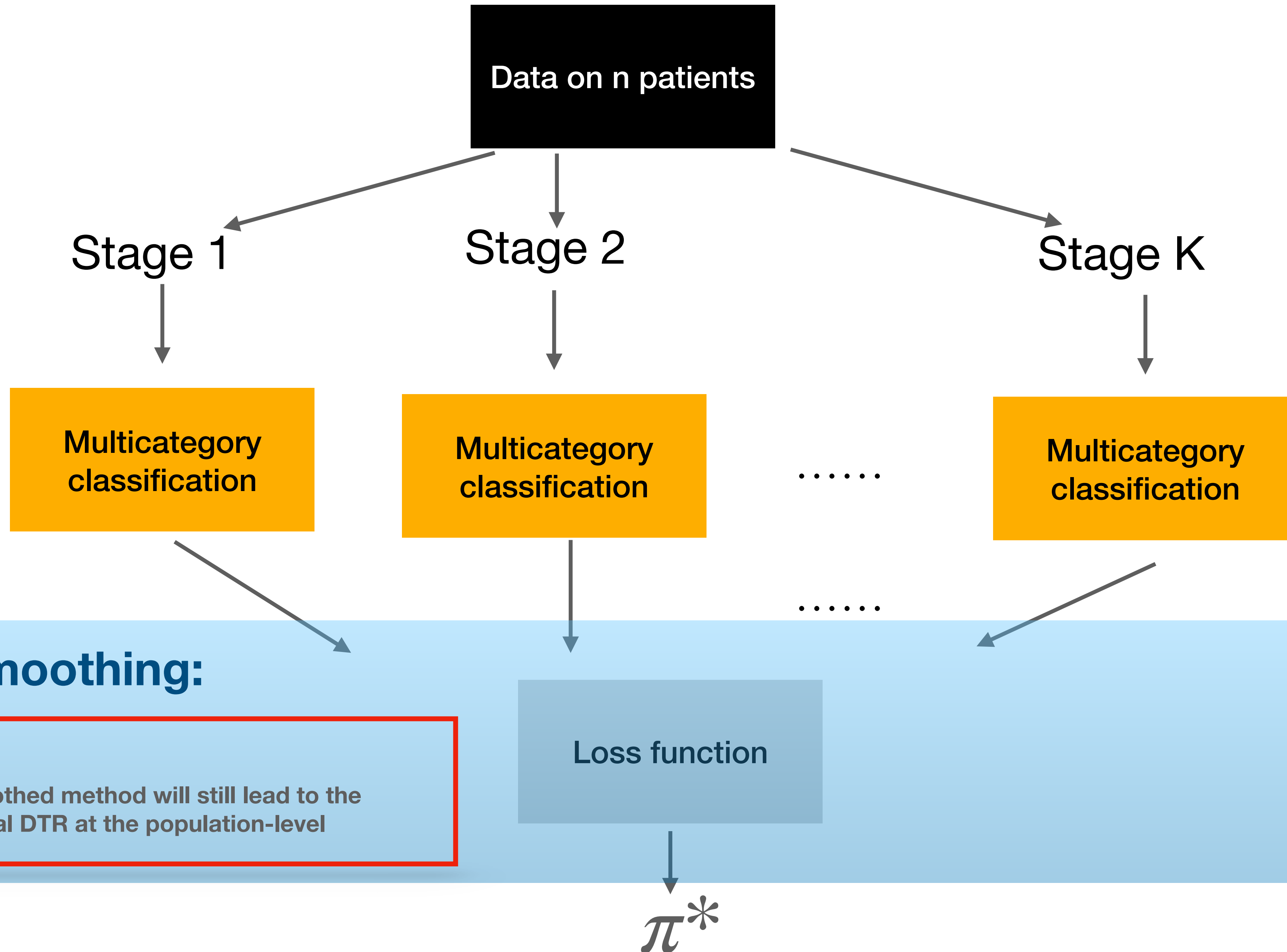
Smoothed loss function



Smoothed loss function



Smoothed loss function



Smoothed loss function

