

Adaptive Security Policies via Belief Aggregation and Rollout

NETCON SEMINAR

Kim Hammar, kimham@kth.se

March 3, 2025

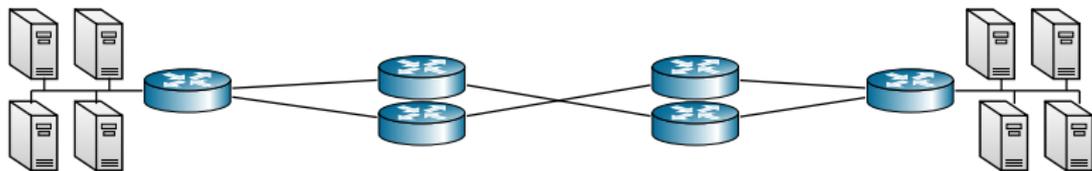
Joint work with Yuchao Li, Tansu Alpcan, Emil Lupu, and Dimitri Bertsekas



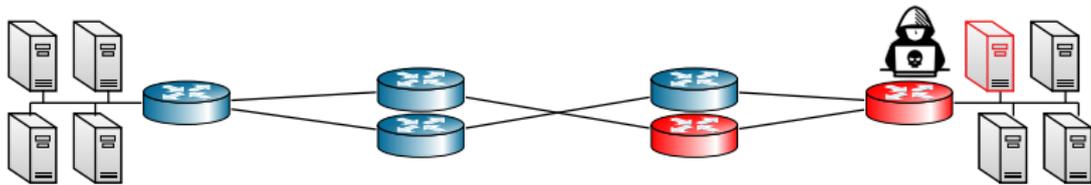
Based on work-in-progress papers ⌚:

Adaptive Security Policies via Belief Aggregation and Rollout (K.H, Y.L, T.A, E.L)
Feature-Based Belief Aggregation for POMDPs (Y.L, K.H, D.B)

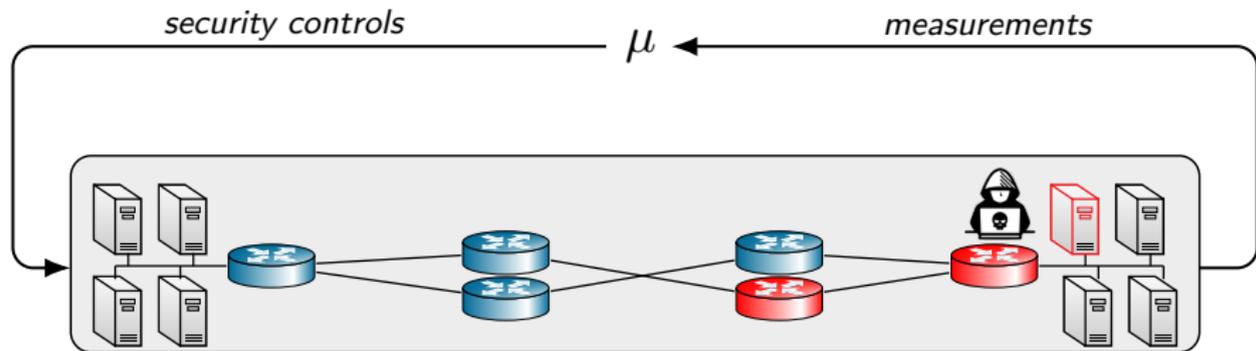
Problem



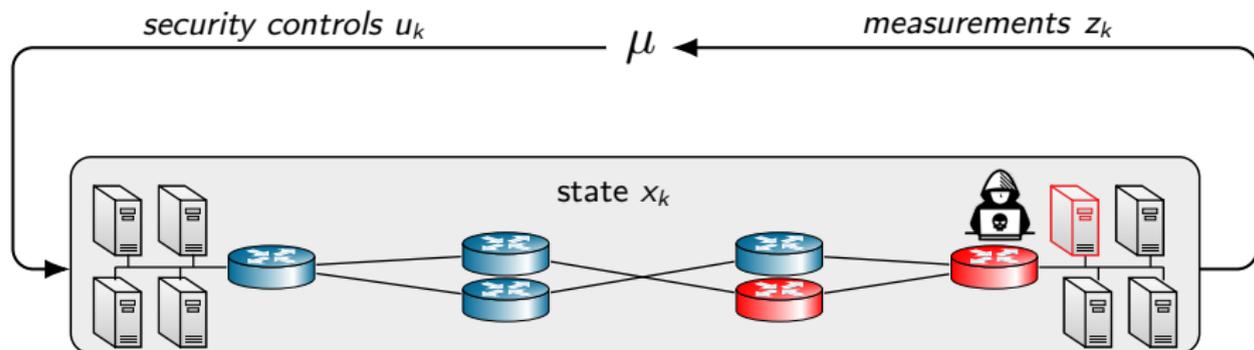
Problem



Problem

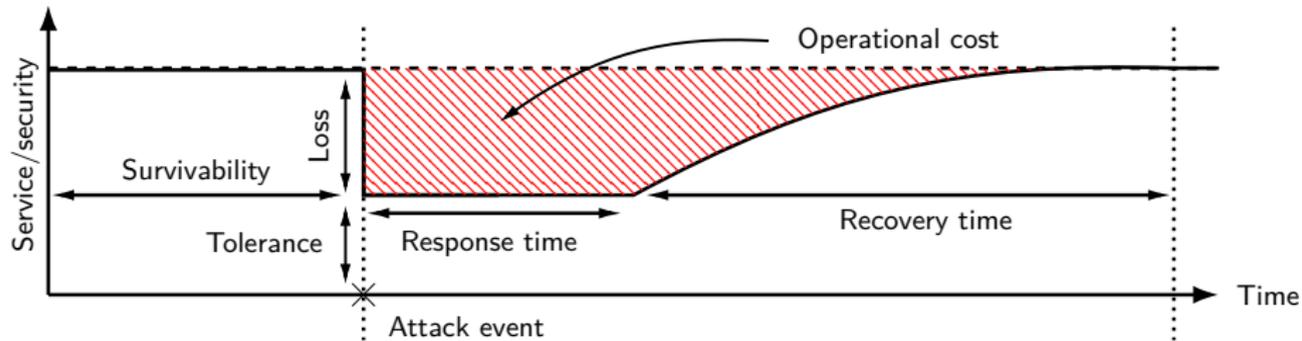


Problem



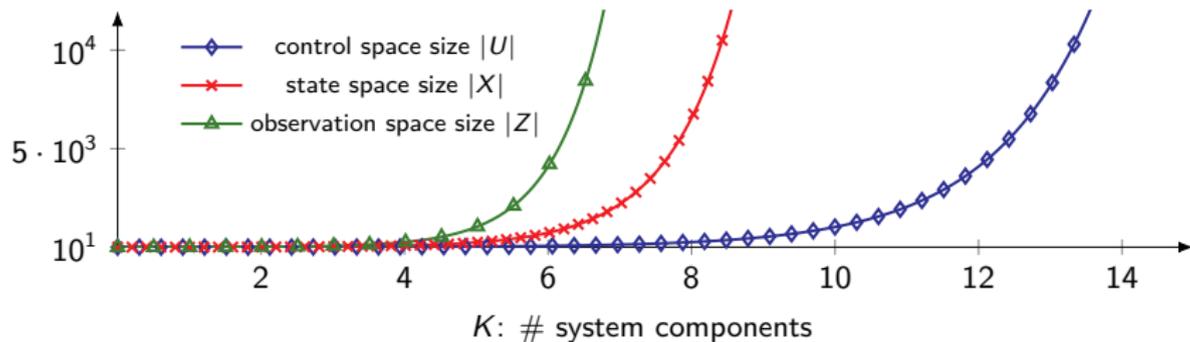
- Finite partially observed Markovian decision problem (POMDP).
- Hidden states $i \in X = \{1, \dots, n\}$, transition probability $p_{ij}(u)$.
- Observation $z \in Z$ is generated with probability $p(z | j, u)$.
- Control $u \in U$.
- **Goal:** find a policy μ that minimizes the discounted cost

$$E \left\{ \sum_{k=0}^{\infty} \alpha^k g(x_k, u_k, x_{k+1}) \right\}.$$



Challenge 1: Curse of Dimensionality

Problem complexity **grows exponentially** with system size.

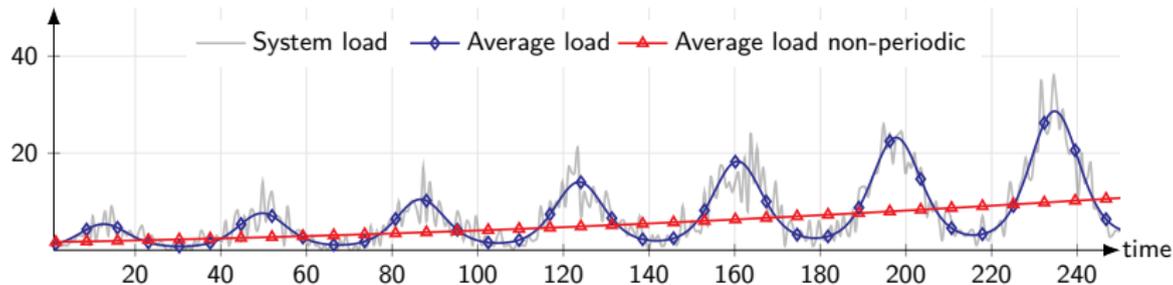


Scalability challenge

Benchmark POMDP (CAGE-2) has over 10^{47} states and 10^{25} observations.

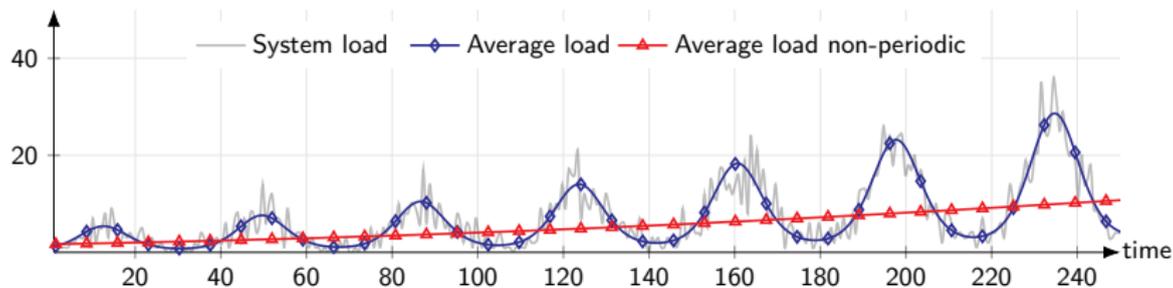
Challenge 2: Changing Dynamics

- Networked systems **change on a regular basis**.
 - ▶ Components fail, bandwidth fluctuates, load patterns shift, software is updated, etc.



Challenge 2: Changing Dynamics

- Networked systems **change on a regular basis**.
 - ▶ Components fail, bandwidth fluctuates, load patterns shift, software is updated, etc.

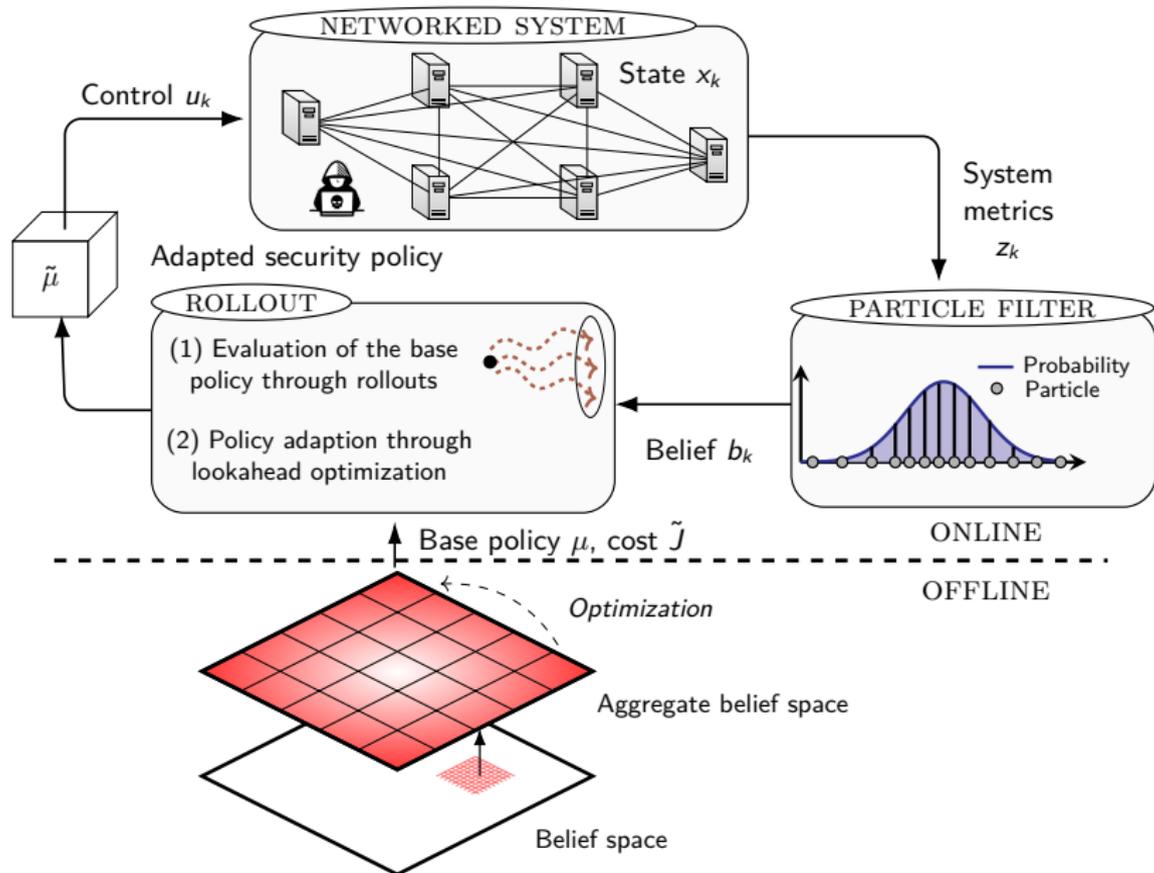


Policy adaption

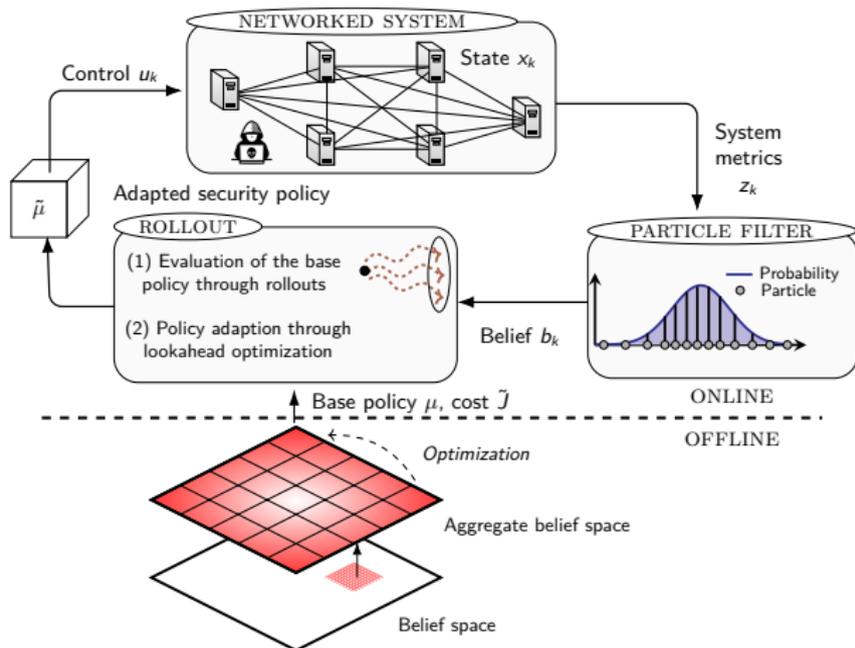
Need an efficient way to adapt the security policy μ when changes occur.

Our Approximation Framework for Large-Scale POMDPs

Our Approximation Framework for Large-Scale POMDPs



Our Approximation Framework for Large-Scale POMDPs



- Achieves **state-of-the-art performance** on the CAGE-2 benchmark.
 - ▶ POMDP with over 10^{47} states and 10^{25} observations.
- Has theoretical **performance guarantees**.
 - ▶ Contrasts with other approximation frameworks, e.g., DEEP RL and LLMs.

Belief state

$b = (b(1), \dots, b(n))$ is a probability distribution over the state space.

Let b be the state, we then obtain **perfect information problem** with **dynamics**

$$b_k = F(b_{k-1}, u_{k-1}, z_k) \quad (\text{Belief estimator})$$

$$\hat{g}(b, u) = \sum_{i=1}^n b(i) \sum_{j=1}^n g(i, u, j) \quad (\text{Cost})$$

$$\hat{p}(z | b, u) = \sum_{i=1}^n b(i) \sum_{j=1}^n p_{ij}(u) p(z | j, u) \quad (\text{Disturbance distribution}).$$

Belief state

$b = (b(1), \dots, b(n)) \in B$ is a probability distribution over the state space.

Let b be the state, we then obtain **perfect information problem** with **dynamics**

$$b_k = F(b_{k-1}, u_{k-1}, z_k) \quad (\text{Belief estimator})$$

$$\hat{g}(b, u) = \sum_{i=1}^n b(i) \sum_{j=1}^n g(i, u, j) \quad (\text{Cost})$$

$$\hat{p}(z | b, u) = \sum_{i=1}^n b(i) \sum_{j=1}^n p_{ij}(u) p(z | j, u) \quad (\text{Disturbance distribution}).$$

Challenges

- Computing an optimal policy is **PSPACE-hard**.

Belief state

$b = (b(1), \dots, b(n)) \in B$ is a probability distribution over the state space.

Let b be the state, we then obtain **perfect information problem** with **dynamics**

$$b_k = F(b_{k-1}, u_{k-1}, z_k) \quad (\text{Belief estimator})$$

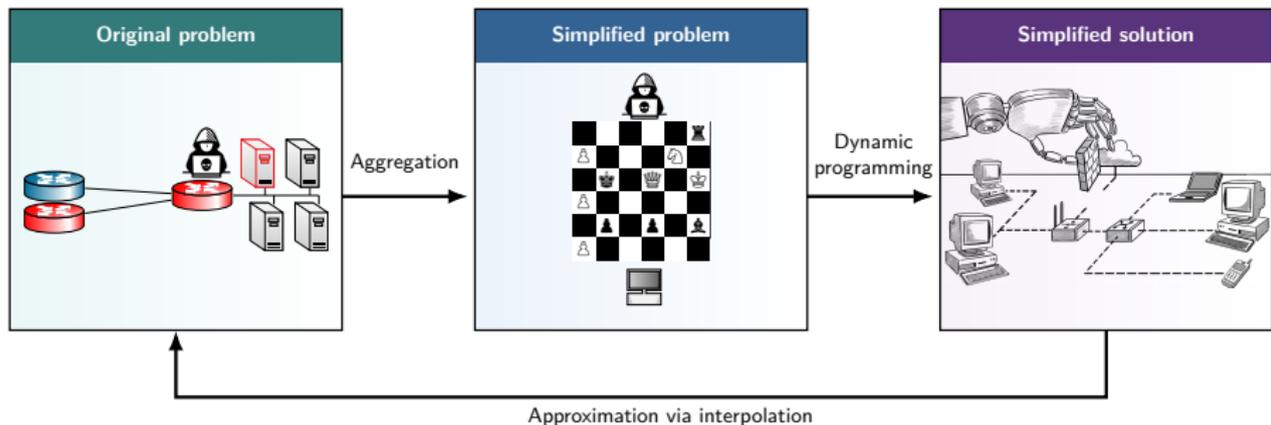
$$\hat{g}(b, u) = \sum_{i=1}^n b(i) \sum_{j=1}^n g(i, u, j) \quad (\text{Cost})$$

$$\hat{p}(z | b, u) = \sum_{i=1}^n b(i) \sum_{j=1}^n p_{ij}(u) p(z | j, u) \quad (\text{Disturbance distribution}).$$

Challenges

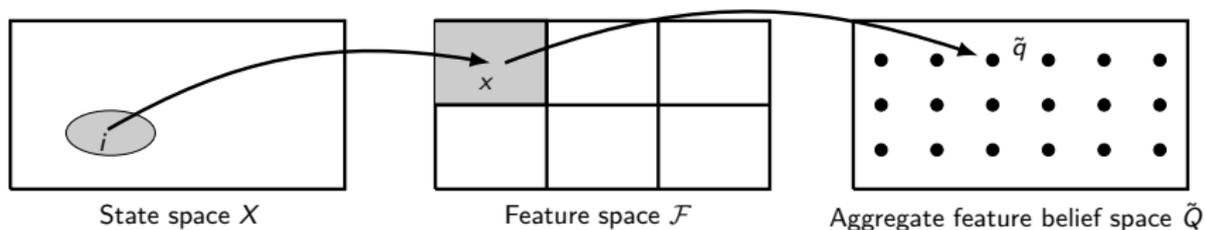
- Computing an optimal policy is **PSPACE-hard**.
- Enumerating the **dimension** of B is intractable ($|X| \geq 10^{47}$ in CAGE-2).

Offline POMDP Approximation via Problem Simplification

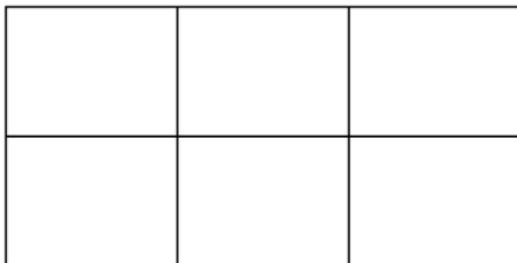


Two-Level Aggregation for Offline Approximation

We use **two-level aggregation** to simplify the POMDP into an **aggregate MDP** with finite state space, which we solve using **dynamic programming**.



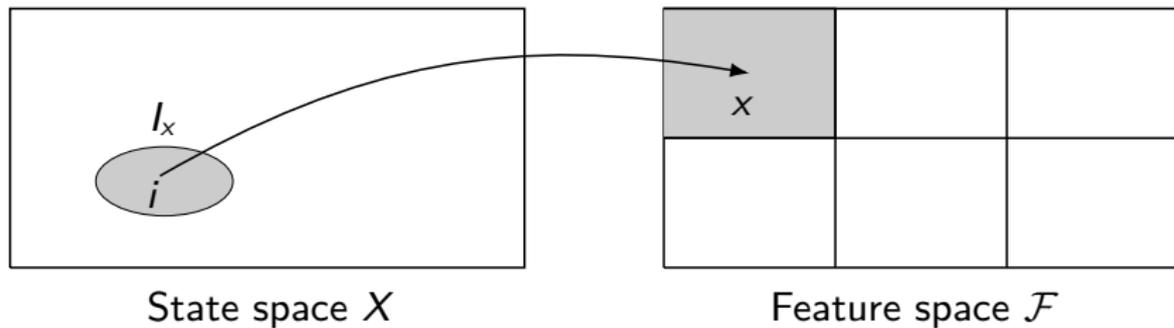
- Two main options to construct the **feature space** \mathcal{F} :
 - ▶ It can be manually designed based on **engineering intuition**.
 - ▶ It can be automatically constructed via a **neural network**.



Feature space \mathcal{F}

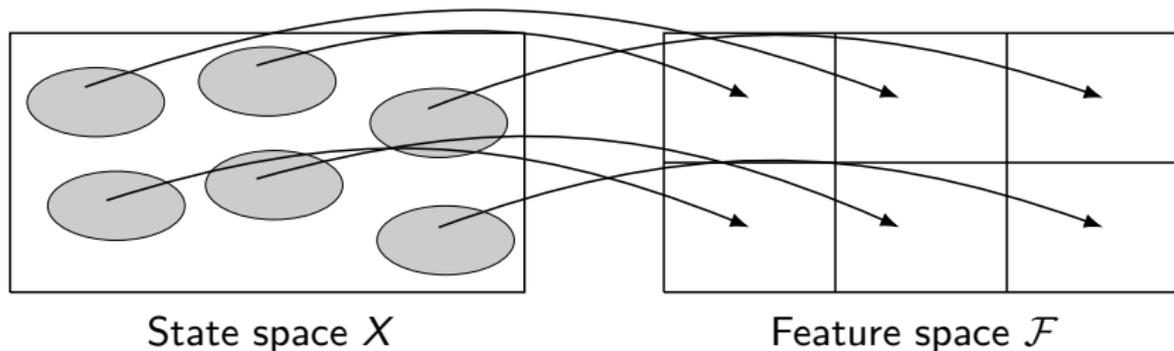
State Aggregation (2/3)

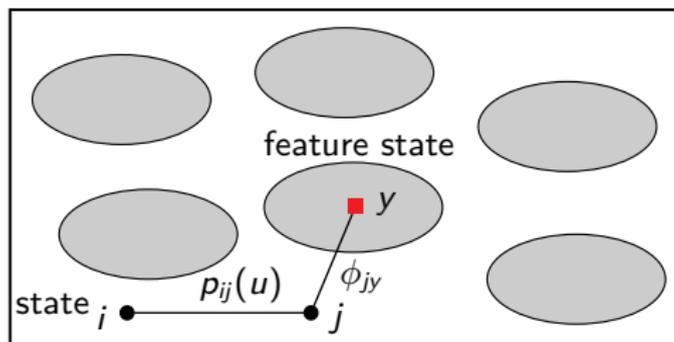
Each **feature state** $x \in \mathcal{F}$ is associated with a **disjoint subset** $I_x \subset X$.



State Aggregation (2/3)

Each **feature state** $x \in \mathcal{F}$ is associated with a **disjoint subset** $I_x \subset X$.





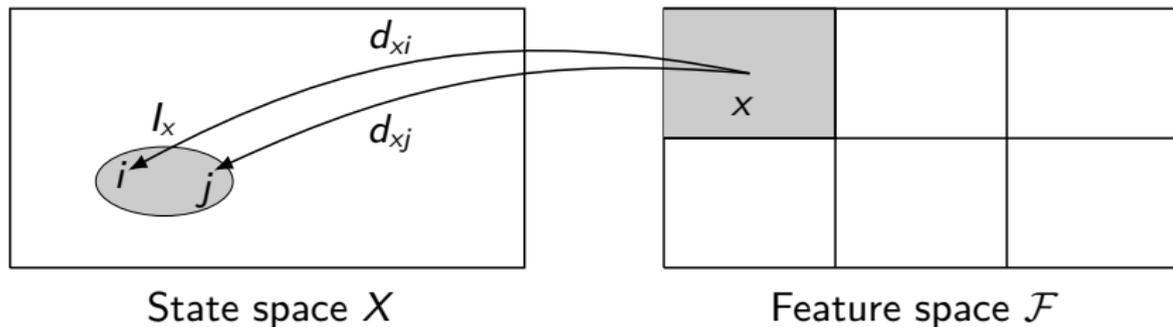
State space X

State aggregation

For each state $j \in X$, we associate

- an **aggregation probability** distribution $\{\phi_{jy} \mid y \in \mathcal{F}\}$, where $\phi_{jy} = 1$ for all $j \in I_y$.

State Aggregation (2/3)



State aggregation

For each state $j \in X$, we associate

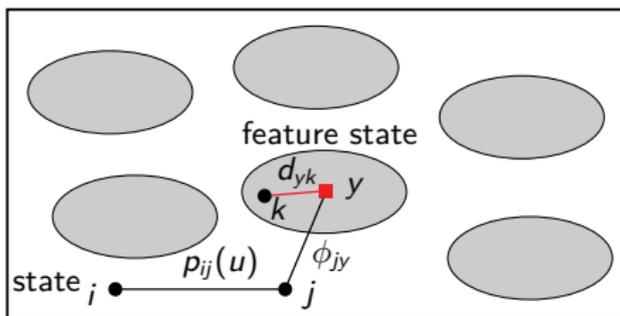
- an **aggregation probability** distribution $\{\phi_{jy} \mid y \in \mathcal{F}\}$, where $\phi_{jy} = 1$ for all $j \in I_y$.

Feature disaggregation

For every feature state $x \in \mathcal{F}$, we associate

- a **disaggregation probability** distribution $\{d_{xi} \mid i \in X\}$, where $d_{xi} = 0$ for all $i \notin I_x$.

State Aggregation (2/3)



State space X

State aggregation

For each state $j \in X$, we associate

- an **aggregation probability** distribution $\{\phi_{jy} \mid y \in \mathcal{F}\}$, where $\phi_{jy} = 1$ for all $j \in I_y$.

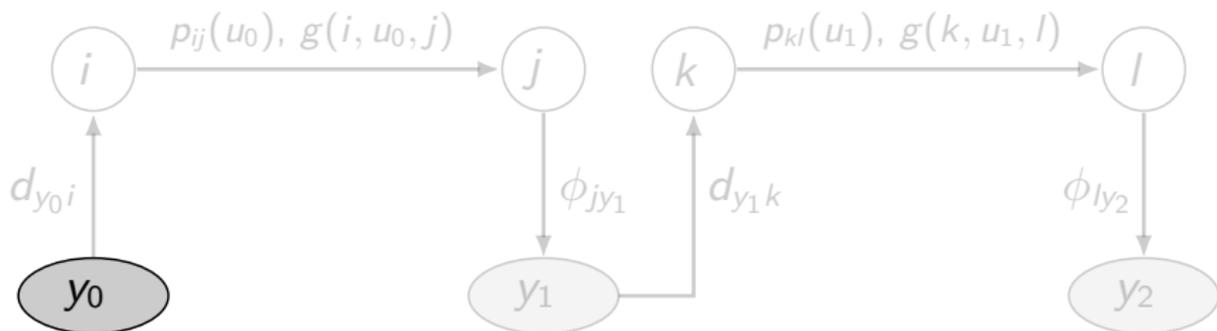
Feature disaggregation

For every *feature state* $x \in \mathcal{F}$, we associate

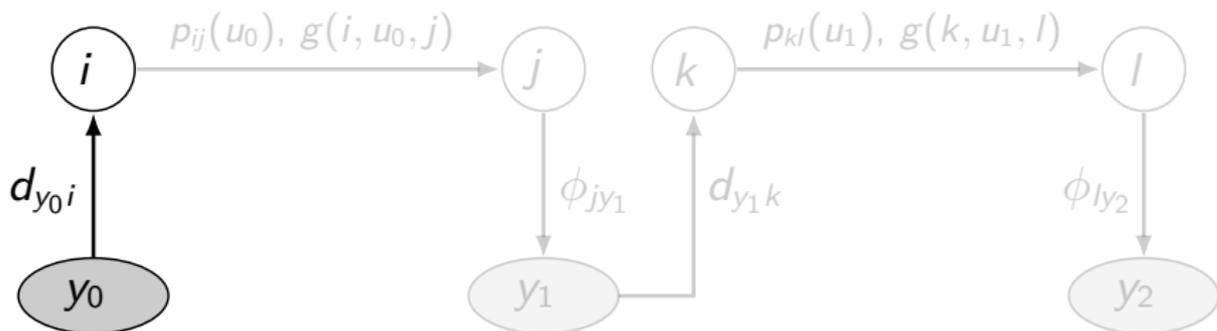
- a **disaggregation probability** distribution $\{d_{xi} \mid i \in X\}$, where $d_{xi} = 0$ for all $i \notin I_x$.

We obtain the **dynamic system**:

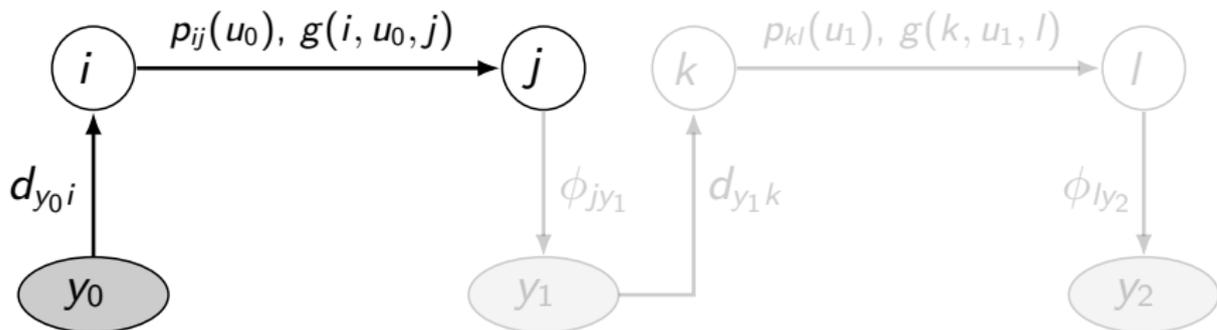




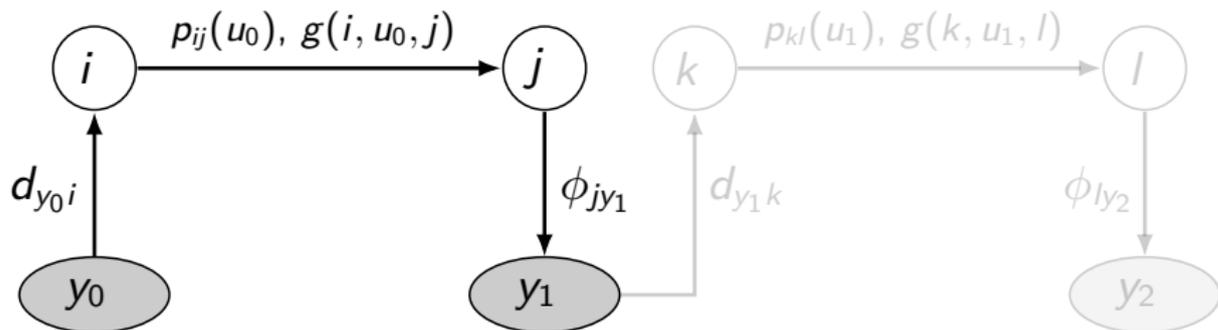
- \implies **Well-defined aggregate problem** with **state space \mathcal{F}** .
 - ▶ We obtain the desired **dimensionality reduction**: $|\mathcal{F}| \ll |\mathcal{X}|$.



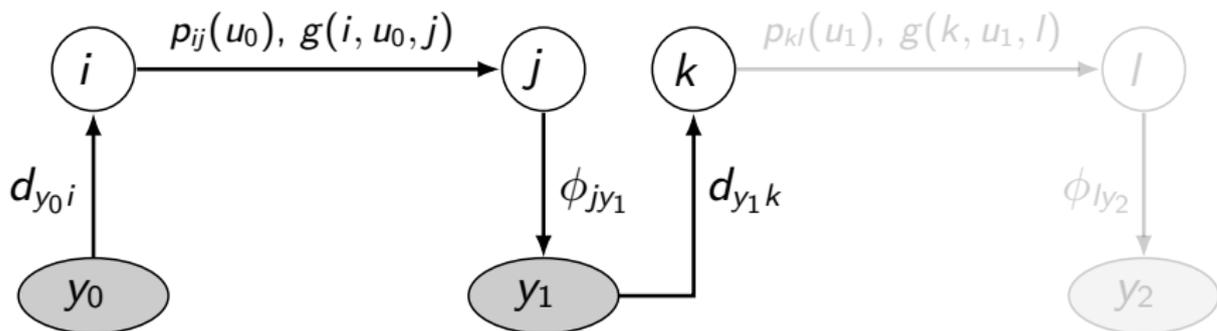
- \implies **Well-defined aggregate problem** with **state space \mathcal{F}** .
 - ▶ We obtain the desired **dimensionality reduction**: $|\mathcal{F}| \ll |\mathcal{X}|$.



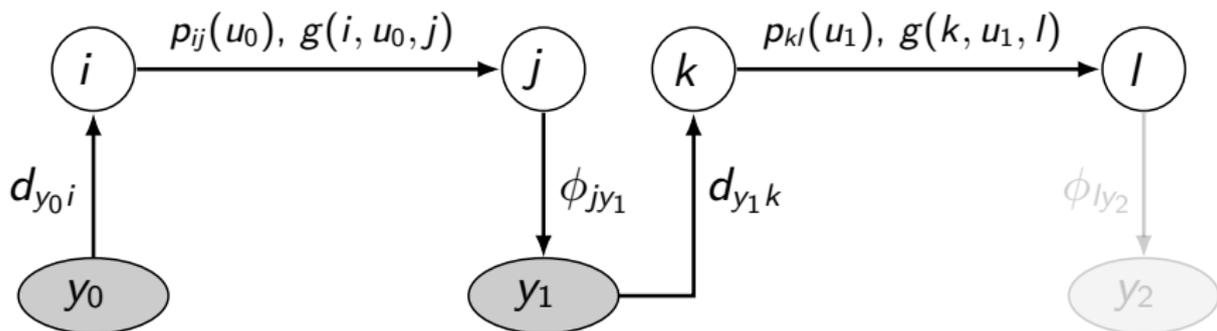
- \implies **Well-defined aggregate problem** with **state space \mathcal{F}** .
 - ▶ We obtain the desired **dimensionality reduction**: $|\mathcal{F}| \ll |\mathcal{X}|$.



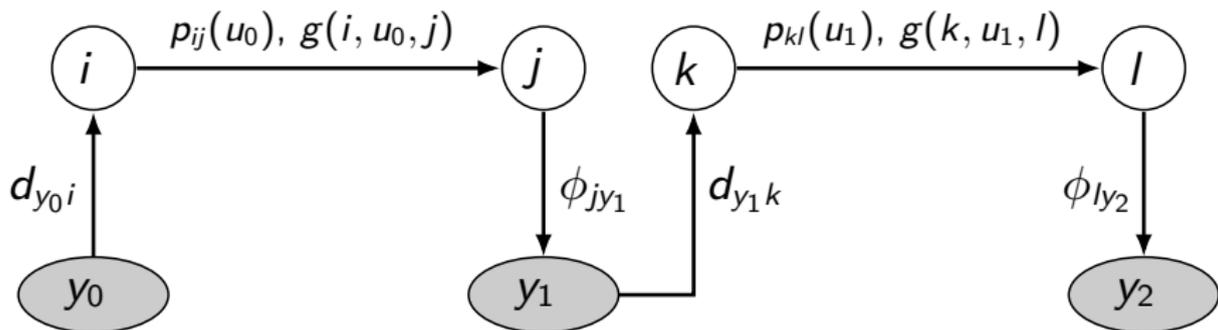
- \implies **Well-defined aggregate problem** with **state space \mathcal{F}** .
 - ▶ We obtain the desired **dimensionality reduction**: $|\mathcal{F}| \ll |\mathcal{X}|$.



- \implies **Well-defined aggregate problem** with **state space \mathcal{F}** .
 - ▶ We obtain the desired **dimensionality reduction**: $|\mathcal{F}| \ll |\mathcal{X}|$.



- \implies **Well-defined aggregate problem** with **state space \mathcal{F}** .
 - ▶ We obtain the desired **dimensionality reduction**: $|\mathcal{F}| \ll |\mathcal{X}|$.



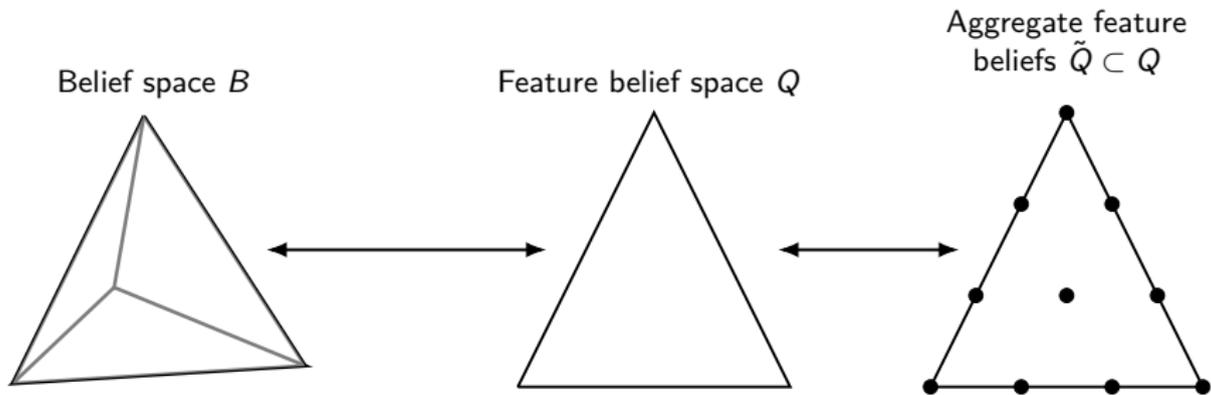
- \implies **Well-defined aggregate problem** with **state space \mathcal{F}** .
 - ▶ We obtain the desired **dimensionality reduction**: $|\mathcal{F}| \ll |\mathcal{X}|$.



How can we lift this **dynamic aggregation system** to the **belief space**?

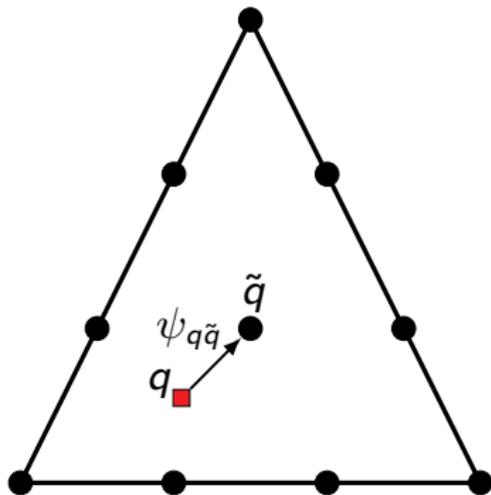
- \implies Well-defined aggregate problem with **state space \mathcal{F}** .

Feature Belief Aggregation (1/3)



Feature Belief Aggregation (2/3)

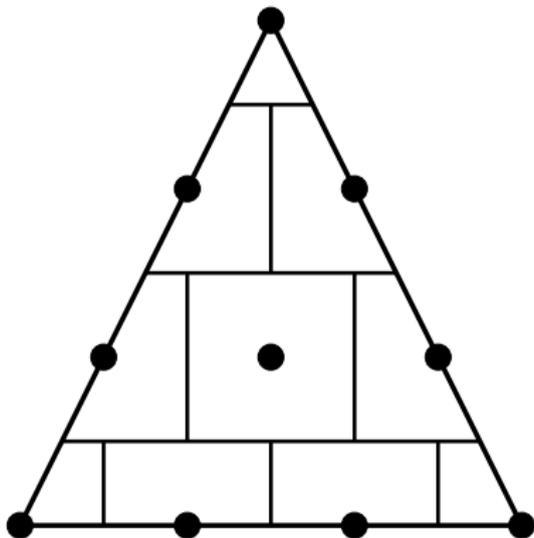
For each **feature belief** $q \in Q$, we associate a *belief aggregation probability* distribution $\{\psi_{q\tilde{q}} \mid \tilde{q} \in \tilde{Q}\}$, where $\psi_{q\tilde{q}} = 1$.



Feature belief space Q

Example (nearest-neighbor aggregation):

$$\psi_{q\tilde{q}} = 1 \iff \tilde{q} = \arg \min_{\tilde{q} \in \tilde{Q}} \|q - \tilde{q}\|.$$

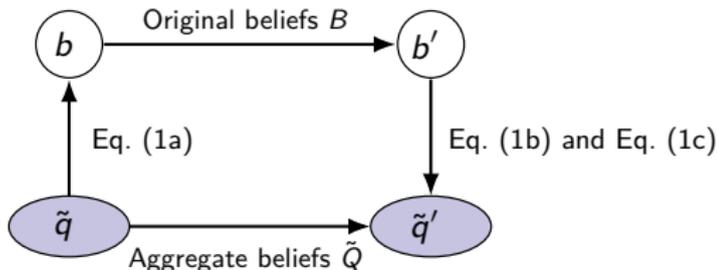
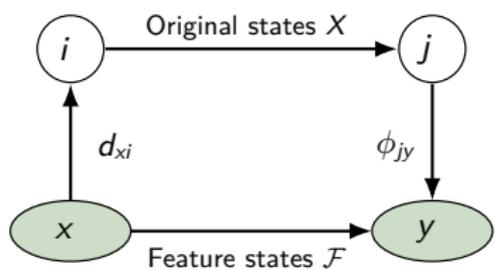


Dynamic belief system:

$$b(i) = \sum_{x \in \mathcal{F}} \tilde{q}(x) d_{xi} \quad \text{for all } i \in X \quad (b \leftarrow q, \tilde{q}) \quad (1a)$$

$$q(y) = \sum_{j=1}^n b(j) \phi_{jy} \quad \text{for all } y \in \mathcal{F} \quad (b \rightarrow q) \quad (1b)$$

$$\tilde{q} \sim \psi_{q\tilde{q}} \quad (q \rightarrow \tilde{q}). \quad (1c)$$



- Let V^* be the **optimal cost-to-go of the aggregate problem**.
- We obtain a cost **approximation of the original POMDP** by

$$\tilde{J}(b) = \sum_{\tilde{q} \in \tilde{Q}} \psi_{\Phi(b)\tilde{q}} V^*(\tilde{q}), \quad (\text{interpolation formula}),$$

where $\Phi : B \mapsto Q$ is defined as

$$\Phi(b)(y) = \sum_{j=1}^n b(j) \phi_{jy} \quad \text{for all } y \in \mathcal{F}.$$

- Let V^* be the **optimal cost-to-go of the aggregate problem**.
- We obtain a cost **approximation of the original POMDP** by

$$\tilde{J}(b) = \sum_{\tilde{q} \in \tilde{Q}} \psi_{\Phi(b)\tilde{q}} V^*(\tilde{q}), \quad (\text{interpolation formula}),$$

where $\Phi : B \mapsto Q$ is defined as

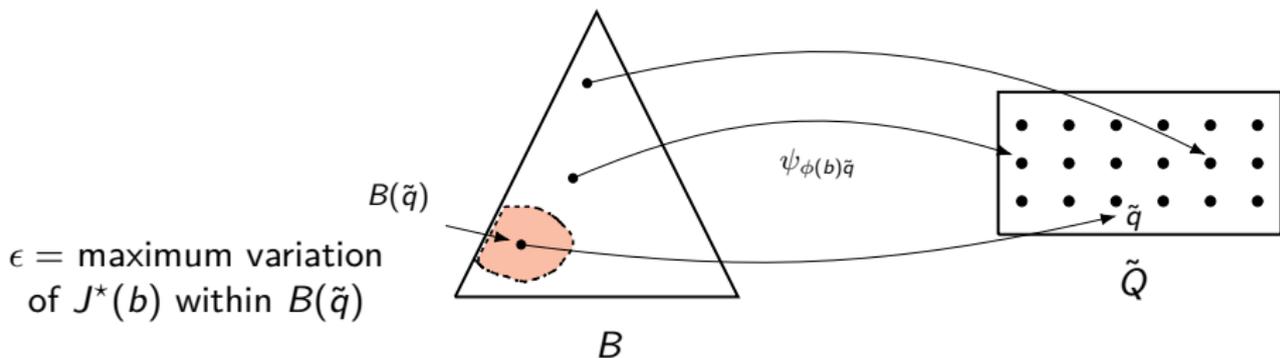
$$\Phi(b)(y) = \sum_{j=1}^n b(j) \phi_{jy} \quad \text{for all } y \in \mathcal{F}.$$

- Similarly, a **base policy** μ for the POMDP can be obtained as

$$\mu(b) \in \arg \min_u E_{b'} \{ \hat{g}(b, u) + \alpha \tilde{J}(b') \} \quad \text{for all } b \in B.$$

- Let V^* be the **optimal cost-to-go of the aggregate problem**.
- We obtain a cost **approximation of the original POMDP** by

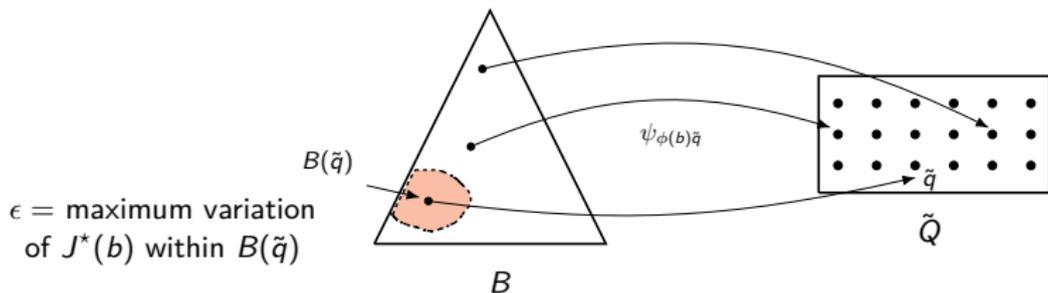
$$\tilde{J}(b) = \sum_{\tilde{q} \in \tilde{Q}} \psi_{\phi(b)\tilde{q}} V^*(\tilde{q}), \quad (\text{interpolation formula}).$$



POMDP Approximation

- Let V^* be the **optimal cost-to-go of the aggregate problem**.
- We obtain a cost **approximation of the original POMDP** by

$$\tilde{J}(b) = \sum_{\tilde{q} \in \tilde{Q}} \psi_{\phi(b)\tilde{q}} V^*(\tilde{q}), \quad (\text{interpolation formula}).$$

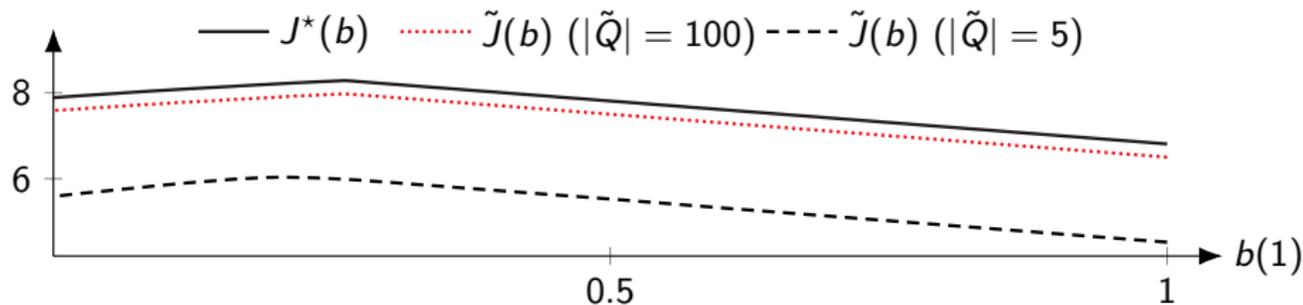


Proposition (Approximation error bound)

Under hard aggregation, the **approximation error** of \tilde{J} is bounded as

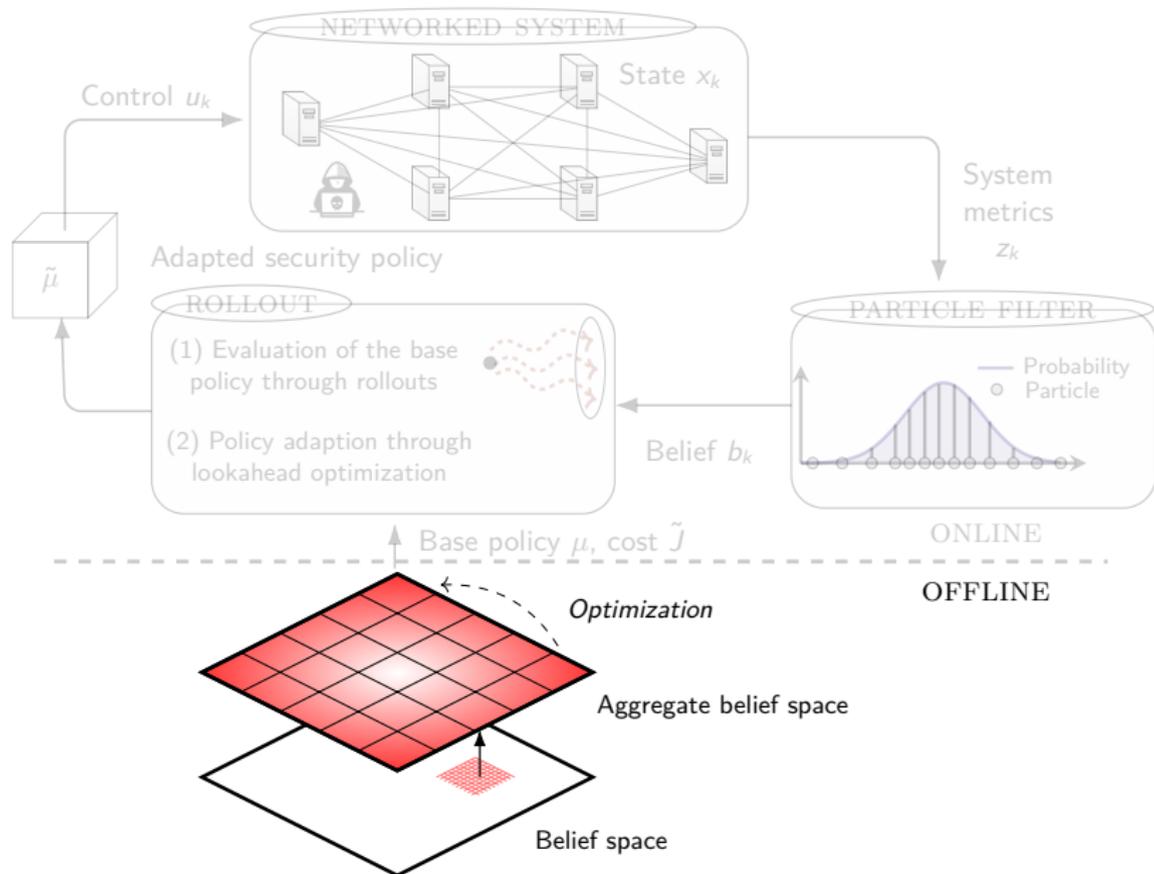
$$|\tilde{J}(b) - J^*(b)| \leq \frac{\epsilon}{1 - \alpha} \quad \forall b \in B(\tilde{q}), \tilde{q} \in \tilde{Q}.$$

Experimental Illustration of the Error Bound

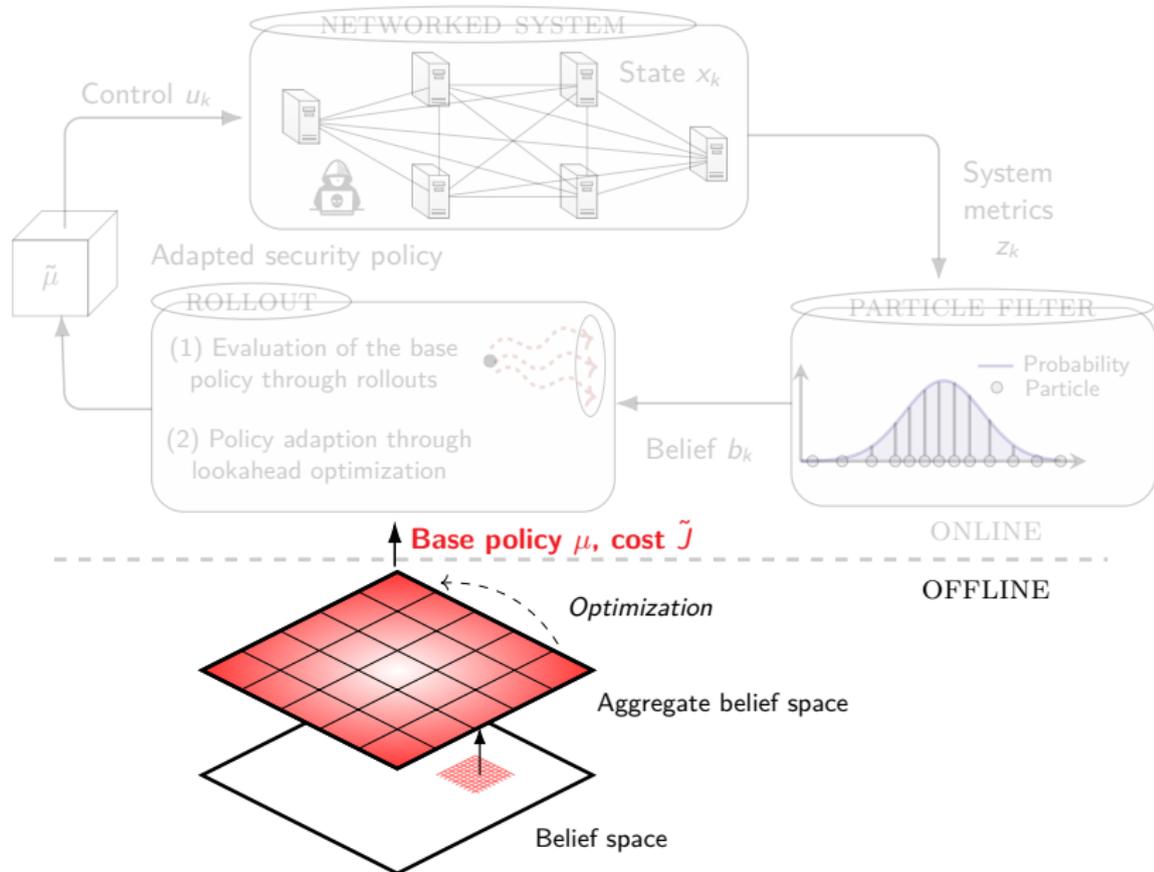


Comparison between the optimal cost-to-go J^* of the POMDP and the approximate cost-to-go \tilde{J} for varying $|\tilde{Q}|$. The numerical results are based on an example POMDP with $|X| = 2$, $\mathcal{F} = X$, \tilde{Q} defined via grid points, and ψ based on the nearest-neighbor mapping.

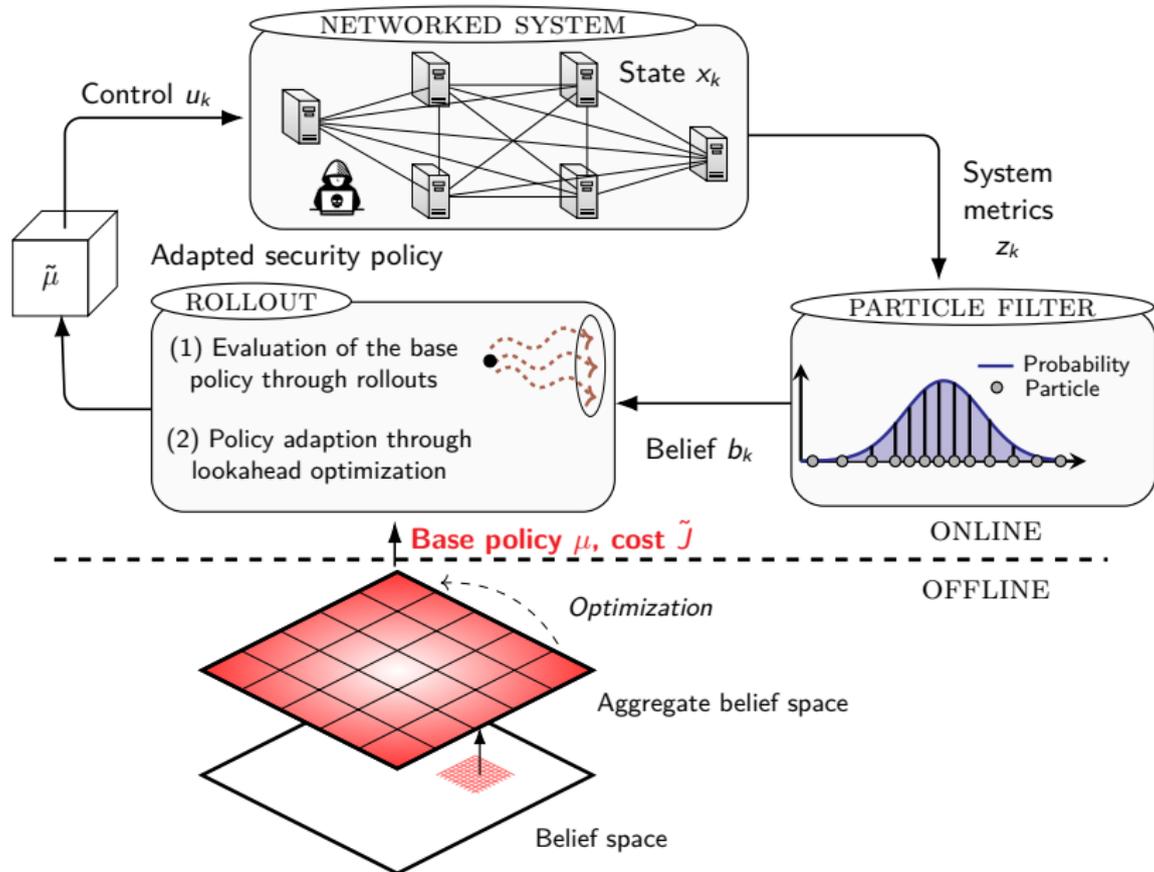
The Big Picture



The Big Picture



The Big Picture



Challenge

Exact computation of the belief b has complexity $O(|X|^2)$, which is **intractable** for realistic systems. (In CAGE-2, $|X| \geq 10^{47}$.)

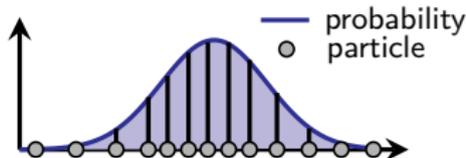
Challenge

Exact computation of the belief b has complexity $O(|X|^2)$, which is **intractable for realistic systems**. (In CAGE-2, $|X| \geq 10^{47}$.)

To manage this complexity, we use a **particle filter** to estimate b as

$$\hat{b}_k(x_k) = \frac{1}{M} \sum_{i=1}^M \mathbb{1}_{x_k = \hat{x}_k^{(i)}}, \quad (2)$$

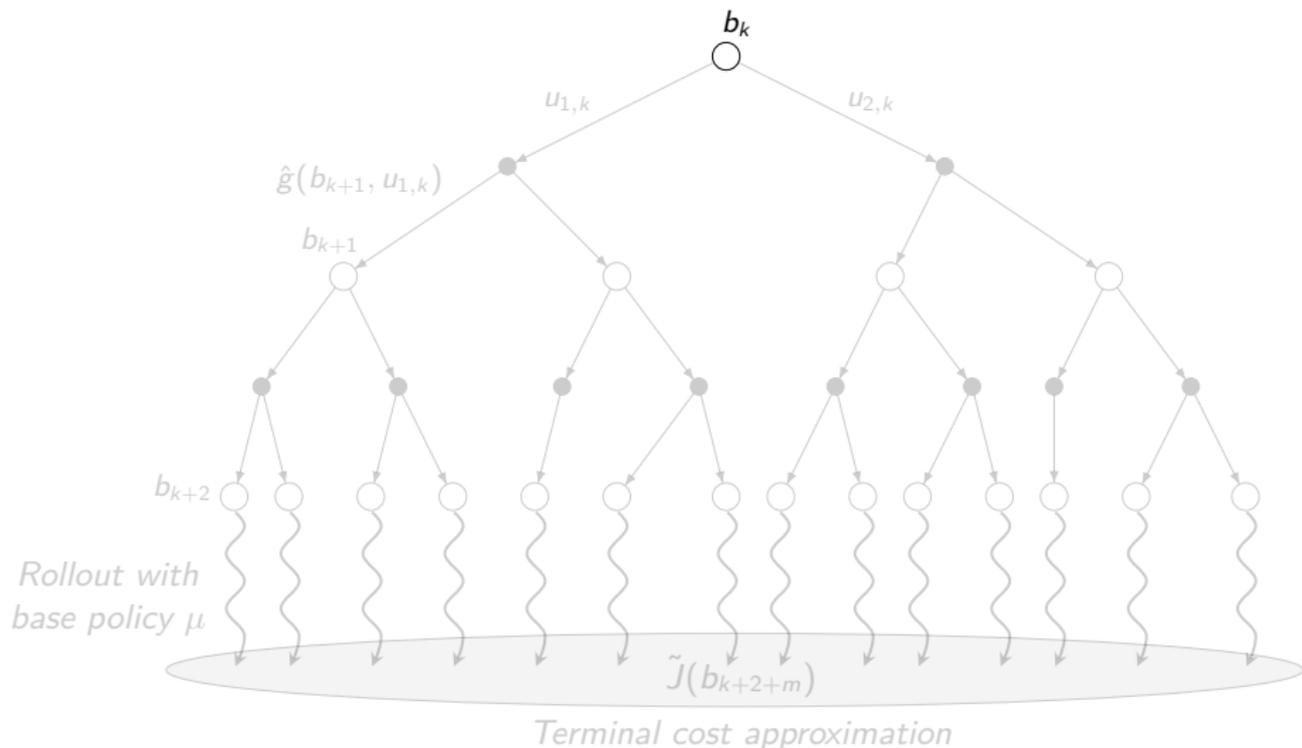
where $\{\hat{x}_k^1, \dots, \hat{x}_k^M\}$ are **particles** sampled with probability proportional to $p(z_k | \hat{x}_k^i, u_{k-1})$.



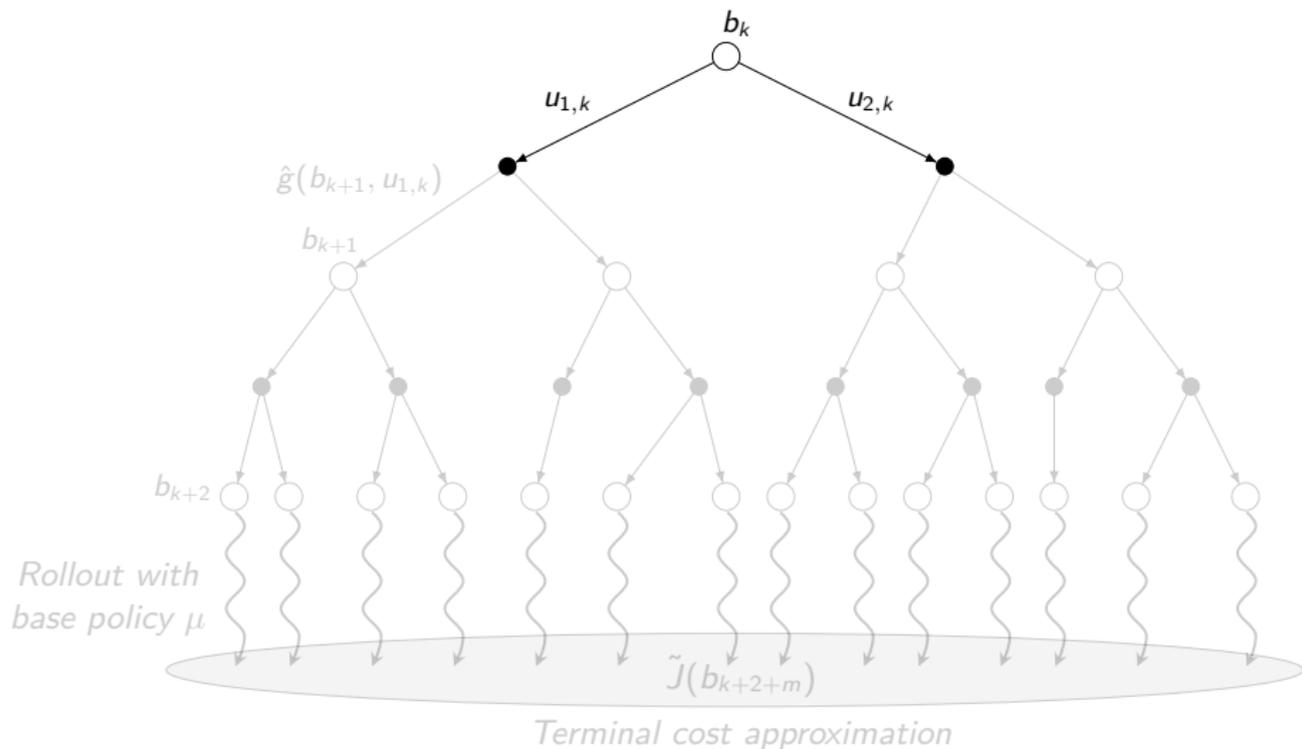
The complexity of Eq. (2) can be adjusted to available compute resources by tuning M . Strong law of large numbers implies asymptotic consistency,

$$\lim_{M \rightarrow \infty} \hat{b} = b.$$

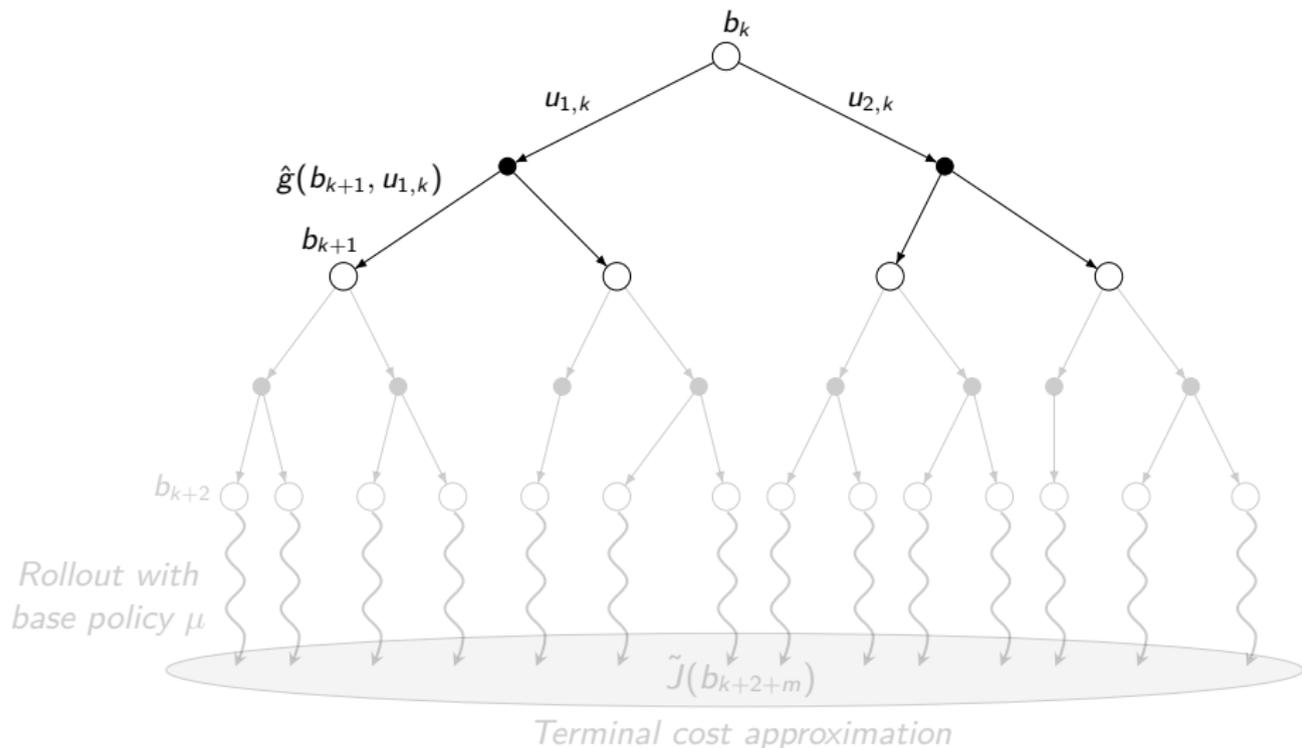
2-STEP LOOKAHEAD OPTIMIZATION WITH ROLLOUT AND TERMINAL COST APPROXIMATION



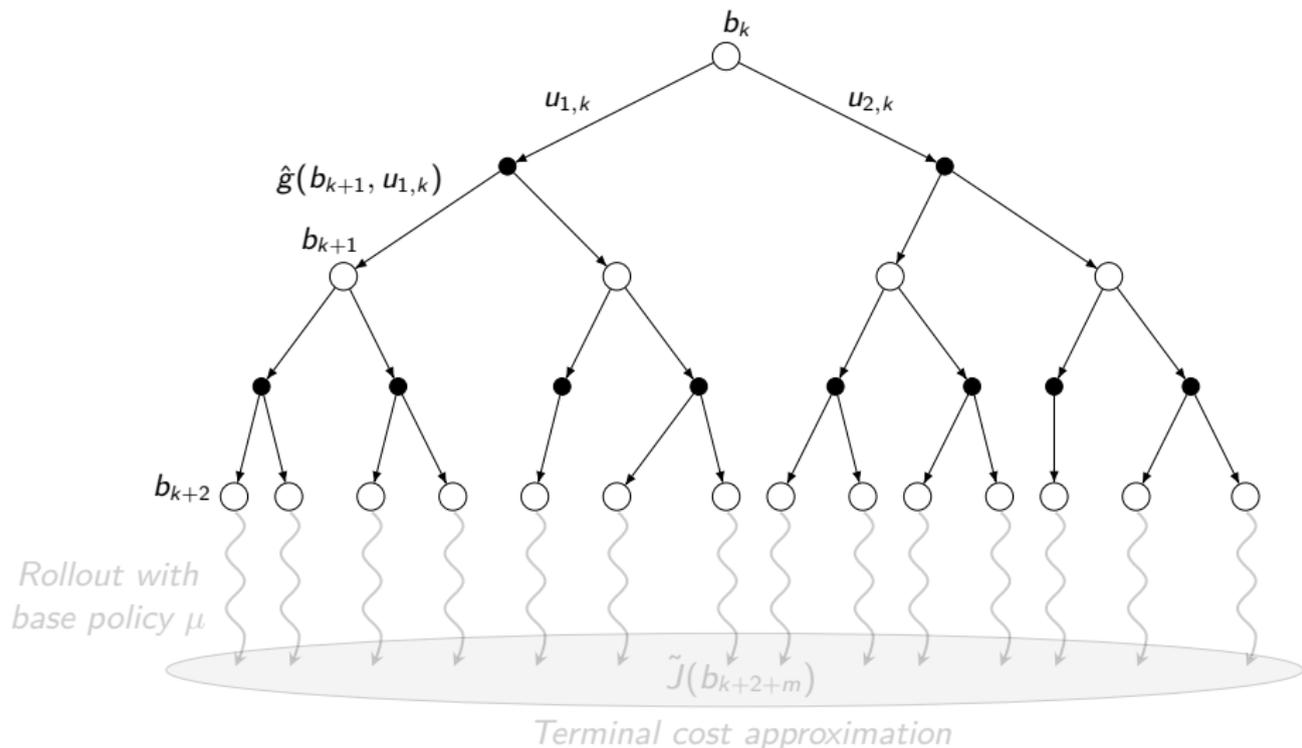
2-STEP LOOKAHEAD OPTIMIZATION WITH ROLLOUT AND TERMINAL COST APPROXIMATION



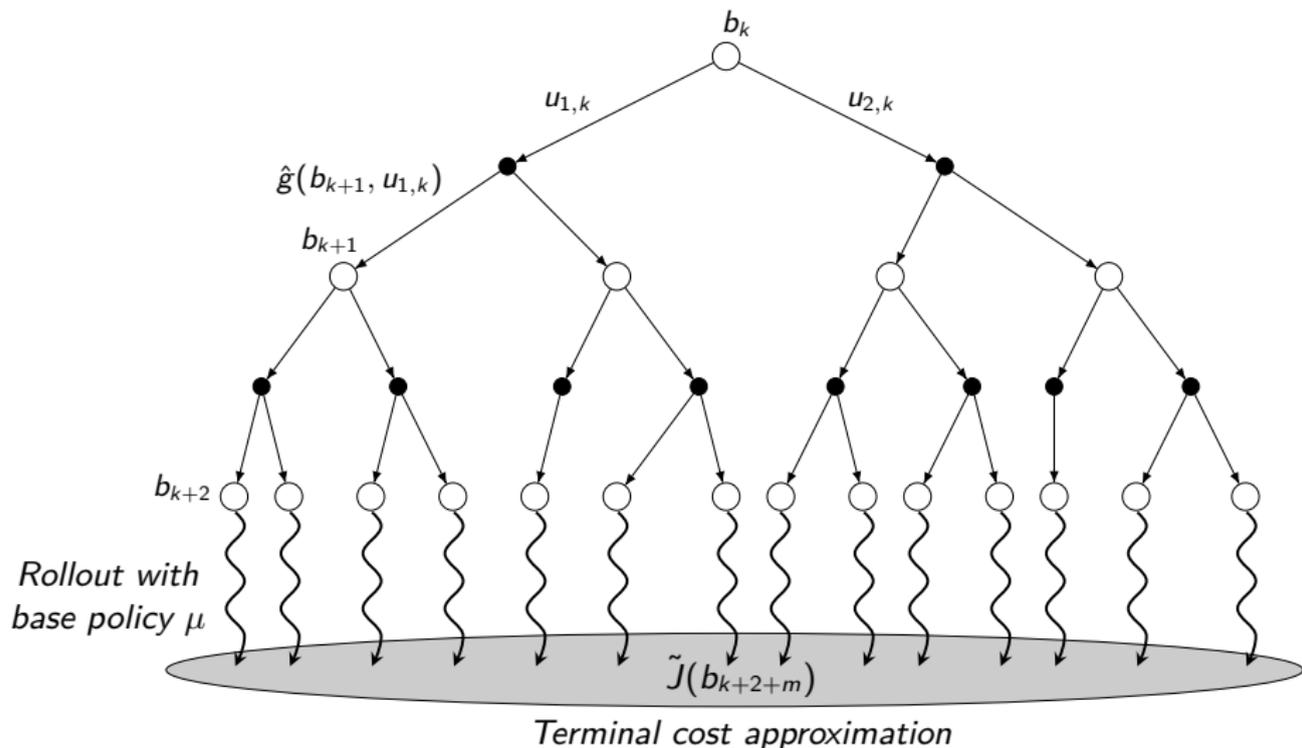
2-STEP LOOKAHEAD OPTIMIZATION WITH ROLLOUT AND TERMINAL COST APPROXIMATION



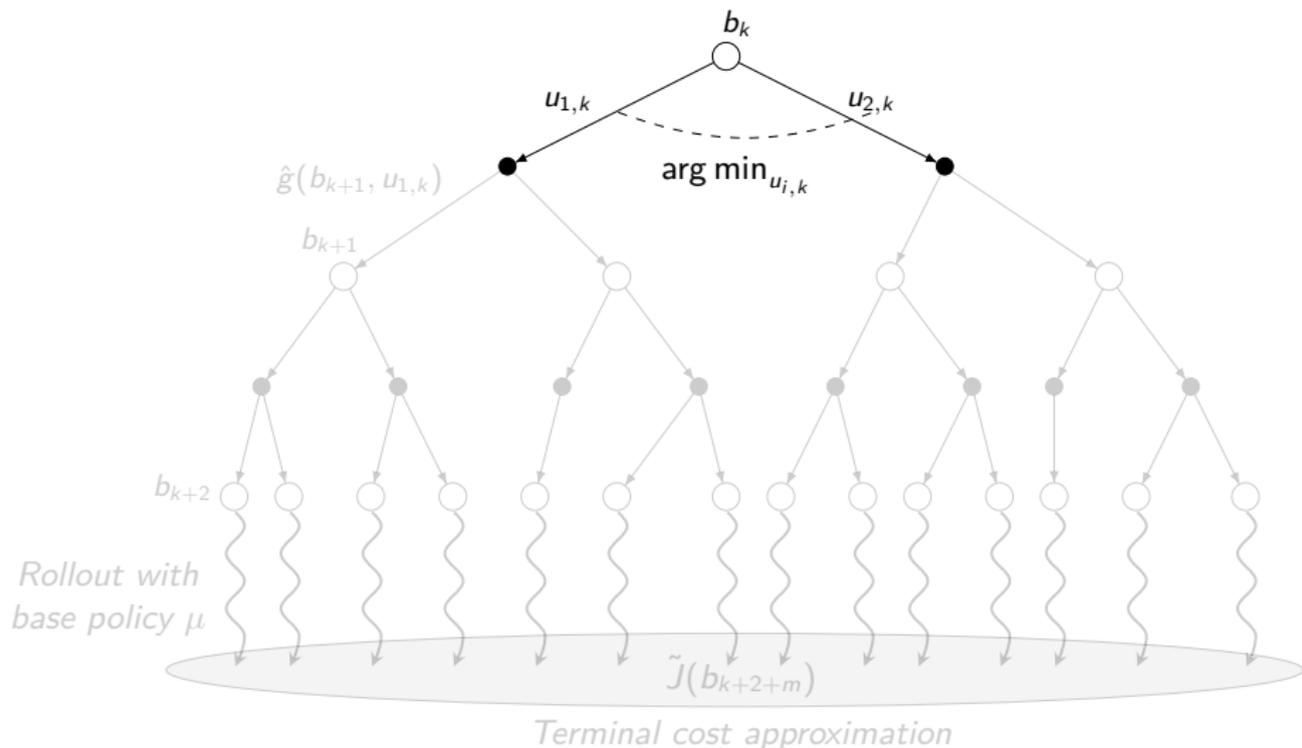
2-STEP LOOKAHEAD OPTIMIZATION WITH ROLLOUT AND TERMINAL COST APPROXIMATION



2-STEP LOOKAHEAD OPTIMIZATION WITH ROLLOUT AND TERMINAL COST APPROXIMATION



2-STEP LOOKAHEAD OPTIMIZATION WITH ROLLOUT AND TERMINAL COST APPROXIMATION



Lookahead optimization

We transform the base policy to an **adapted rollout policy** $\tilde{\mu}$ via ℓ -step lookahead

$$\tilde{\mu}(b_k) \in \arg \min_{u_k, \mu_{k+1}, \dots, \mu_{k+\ell-1}} E_{z_{k+1}, \dots, z_{k+\ell}} \left\{ \hat{g}(b_k, u_k) + \sum_{j=k+1}^{t+\ell-1} \alpha^{j-k} \hat{g}(b_j, \mu_j(b_j)) + \alpha^\ell \tilde{J}_\mu(b_{k+\ell}) \right\}.$$

Rollout

The cost-to-go in the lookahead minimization is **estimated via m -step rollout** with the **base policy** μ and **terminal cost approximation** \tilde{J} :

$$\tilde{J}_\mu(b_k) = \frac{1}{L} \sum_{j=1}^L \sum_{l=k}^{k+m} \alpha^{l-k} \hat{g}(b_l^j, \mu(b_l^j)) + \alpha^{m-k} \tilde{J}(b_{k+m}^j).$$

Proposition (Bertsekas, 2019)

- 1 If the rollout policy evaluation is exact, i.e., if $\tilde{J}_\mu = J_\mu$, then the rollout policy is *guaranteed to improve the base policy*.
- 2 The **sub-optimality of the rollout policy** $\tilde{\mu}$ is bounded as

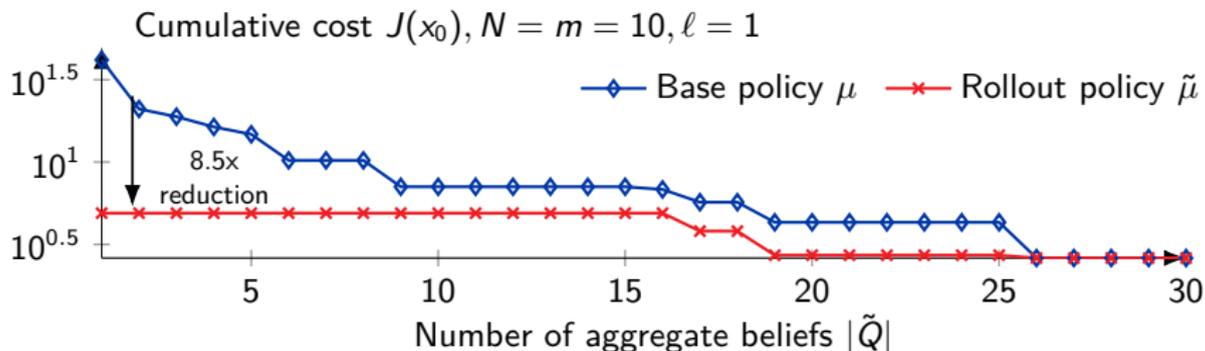
$$\|J_{\tilde{\mu}} - J^*\| \leq \frac{2\alpha^\ell}{1-\alpha} \|\tilde{J}_\mu - J^*\|.$$

Rollout Policy Improvement

Proposition (Bertsekas, 2019)

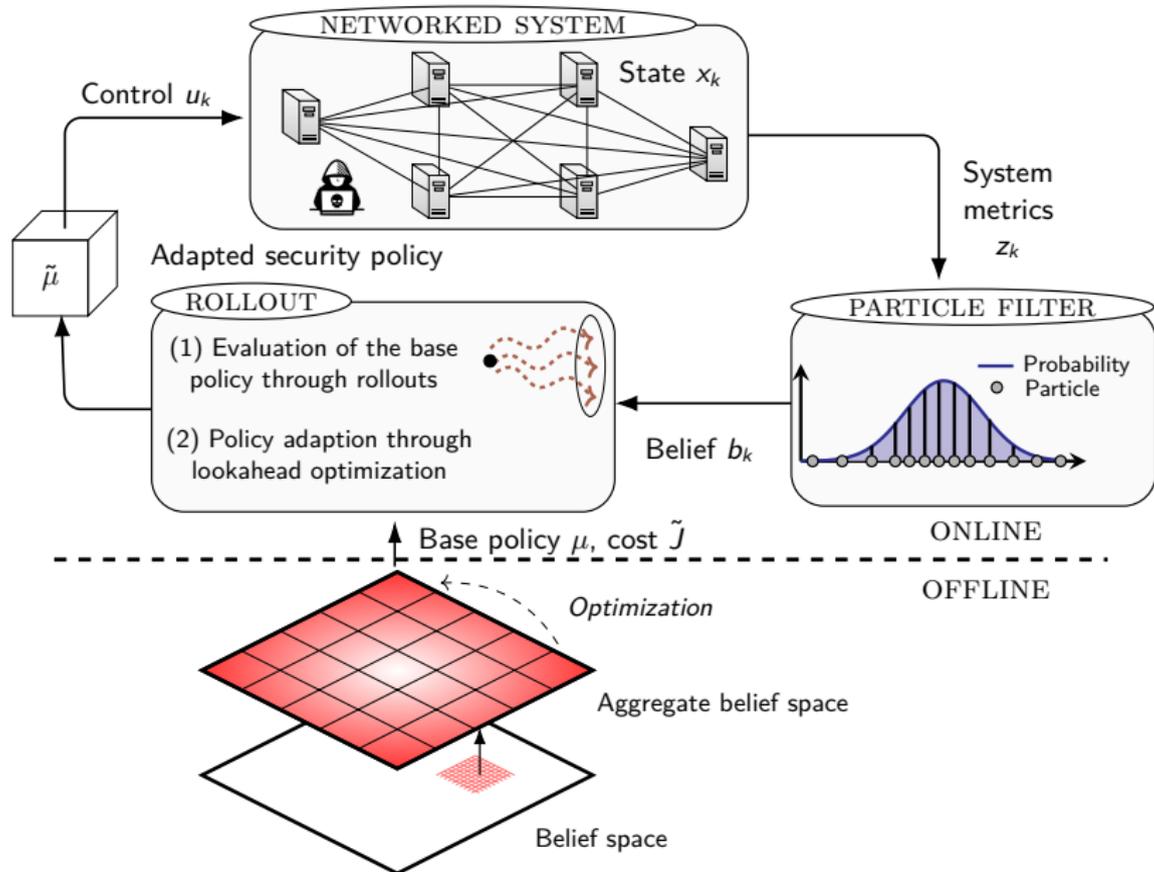
- 1 If the rollout policy evaluation is exact, i.e., if $\tilde{J}_\mu = J_\mu$, then the rollout policy is *guaranteed to improve the base policy*.
- 2 The **sub-optimality of the rollout policy $\tilde{\mu}$ is bounded** as

$$\|J_{\tilde{\mu}} - J^*\| \leq \frac{2\alpha^\ell}{1-\alpha} \|\tilde{J}_\mu - J^*\|.$$



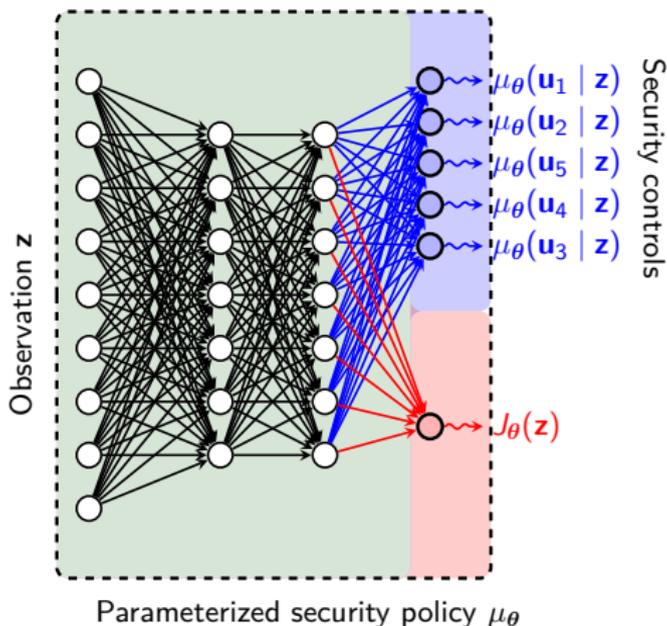
Performance of rollout for an example POMDP.

Framework Summary



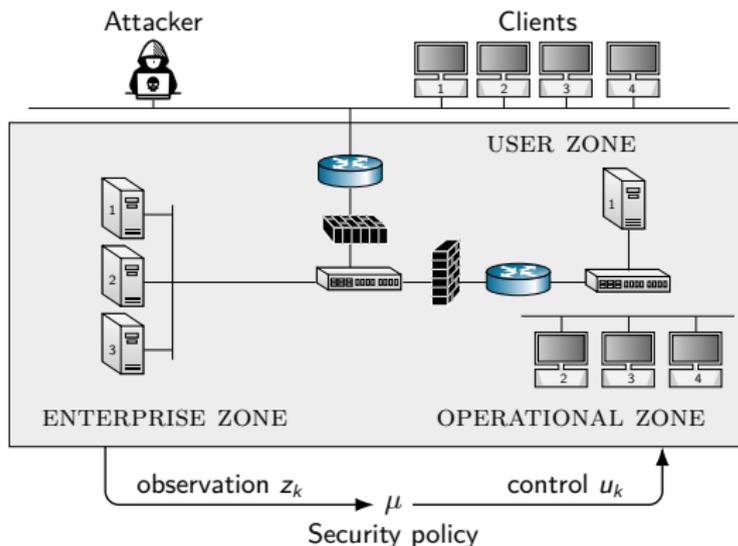
Experimental Evaluation Against the CAGE-2 Benchmark

- **Standard benchmark** for comparing methods: **CAGE-2**.
 - ▶ *Problem*: find an effective security policy to network intrusions.
 - ▶ POMDP with over 10^{47} states and 10^{25} observations.
 - ▶ **Leaderboard** with more than 35 different methods.
- **Current state-of-the-art**: deep reinforcement learning (variants of PPO).



Instantiation of Our Framework for CAGE-2

- We define \tilde{Q} based on intuition, where $|\tilde{Q}| = 427500$.
- We define ψ based on the **nearest-neighbor mapping**.
- We solve the aggregate problem using **value iteration**.
 - ▶ This gives us the base policy μ and cost \tilde{J} .
- We use $\ell = 2$ **lookahead steps** and $M = 50$ **particles**.



The CAGE-2 system.

Experimental Results (1/2)

<i>Method</i>	<i>Offline/Online compute (min/s)</i>	<i>State estimation</i>	<i>Cost</i>
μ	8.5/0.01	PARTICLE FILTER	15.19 \pm 0.82
PPO	1000/0.01	LATEST OBSERVATION	280 \pm 114
PPO	1000/0.01	PARTICLE FILTER	119 \pm 58
PPG	1000/0.01	LATEST OBSERVATION	338 \pm 147
PPG	1000/0.01	PARTICLE FILTER	299 \pm 108
DQN	1000/0.01	LATEST OBSERVATION	479 \pm 267
DQN	1000/0.01	PARTICLE FILTER	462 \pm 244
CARDIFF	300/0.01	LATEST OBSERVATION	13.69 \pm 0.53
CARDIFF	300/0.01	PARTICLE FILTER	13.31 \pm 0.87
POMCP	0/15	PARTICLE FILTER	30.88 \pm 1.41
POMCP	0/30	PARTICLE FILTER	29.51 \pm 2.00
OURS ($m = 0$)	8.5/0.95	PARTICLE FILTER	13.24 \pm 0.57
OURS ($m = 10$)	8.5/8.29	PARTICLE FILTER	13.23 \pm 0.62
OURS ($m = 20$)	8.5/14.80	PARTICLE FILTER	13.23 \pm 0.57

Numbers indicate the mean and the standard deviation from 1000 evaluations.

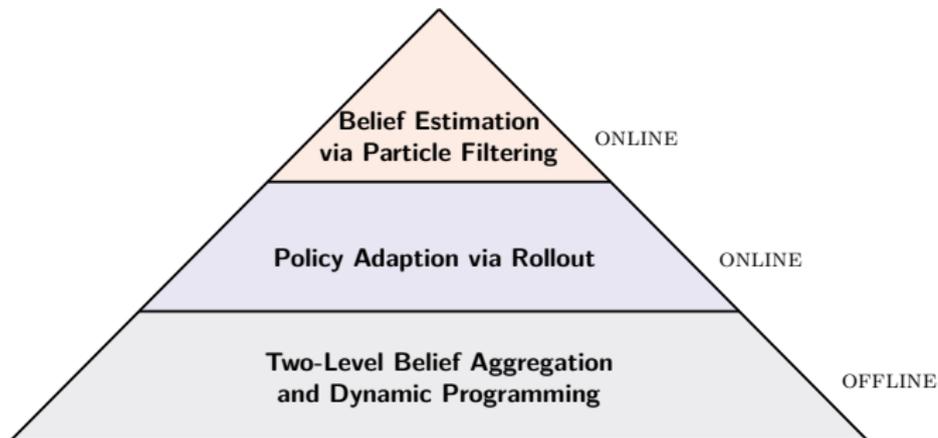
Experimental Results (2/2)

<i>Method</i>	<i>Offline/Online compute (min/s)</i>	<i>State estimation</i>	<i>Cost</i>
μ	8.5/0.01	PARTICLE FILTER	61.72 \pm 3.96
PPO	1000/0.01	LATEST OBSERVATION	341 \pm 133
PPO	1000/0.01	PARTICLE FILTER	326 \pm 116
PPG	1000/0.01	LATEST OBSERVATION	328 \pm 178
PPG	1000/0.01	PARTICLE FILTER	312 \pm 163
DQN	1000/0.01	LATEST OBSERVATION	516 \pm 291
DQN	1000/0.01	PARTICLE FILTER	492 \pm 204
CARDIFF	300/0.01	LATEST OBSERVATION	57.45 \pm 2.44
CARDIFF	300/0.01	PARTICLE FILTER	56.45 \pm 2.81
POMCP	0/15	PARTICLE FILTER	53.08 \pm 3.78
POMCP	0/30	PARTICLE FILTER	53.18 \pm 3.42
OURS ($m = 0$)	8.5/0.95	PARTICLE FILTER	51.87 \pm 1.42
OURS ($m = 10$)	8.5/8.29	PARTICLE FILTER	38.81 \pm 1.68
OURS ($m = 20$)	8.5/14.80	PARTICLE FILTER	37.89 \pm 1.54

Numbers indicate the mean and the standard deviation from 1000 evaluations.

Conclusion

- We present a **scalable framework** for computing **adaptive security policies**, which has formal **performance guarantees** and achieves **state-of-the-art performance**.
- It consists of three components:



Work in progress! ⌚

Theoretical and experimental details will be available in preprints soon.

Source code is available at:

- https://github.com/Limmen/rollout_aggregation; and
- <https://github.com/Limmen/csle>