# Linear Complementarity for Regularized Policy Evaluation and Improvement

*Jeff Johns*

*Christopher Painter-Wakefield*

*Ronald Parr*

Duke
UNIVERSITY

# This Talk

**What** — Solving multiple, related $L_1$ regularization problems

**Why** — Efficient policy iteration algorithm for Markov decision processes (MDPs)

**How** — Formulate as a linear complementarity problem (LCP), use warm starts
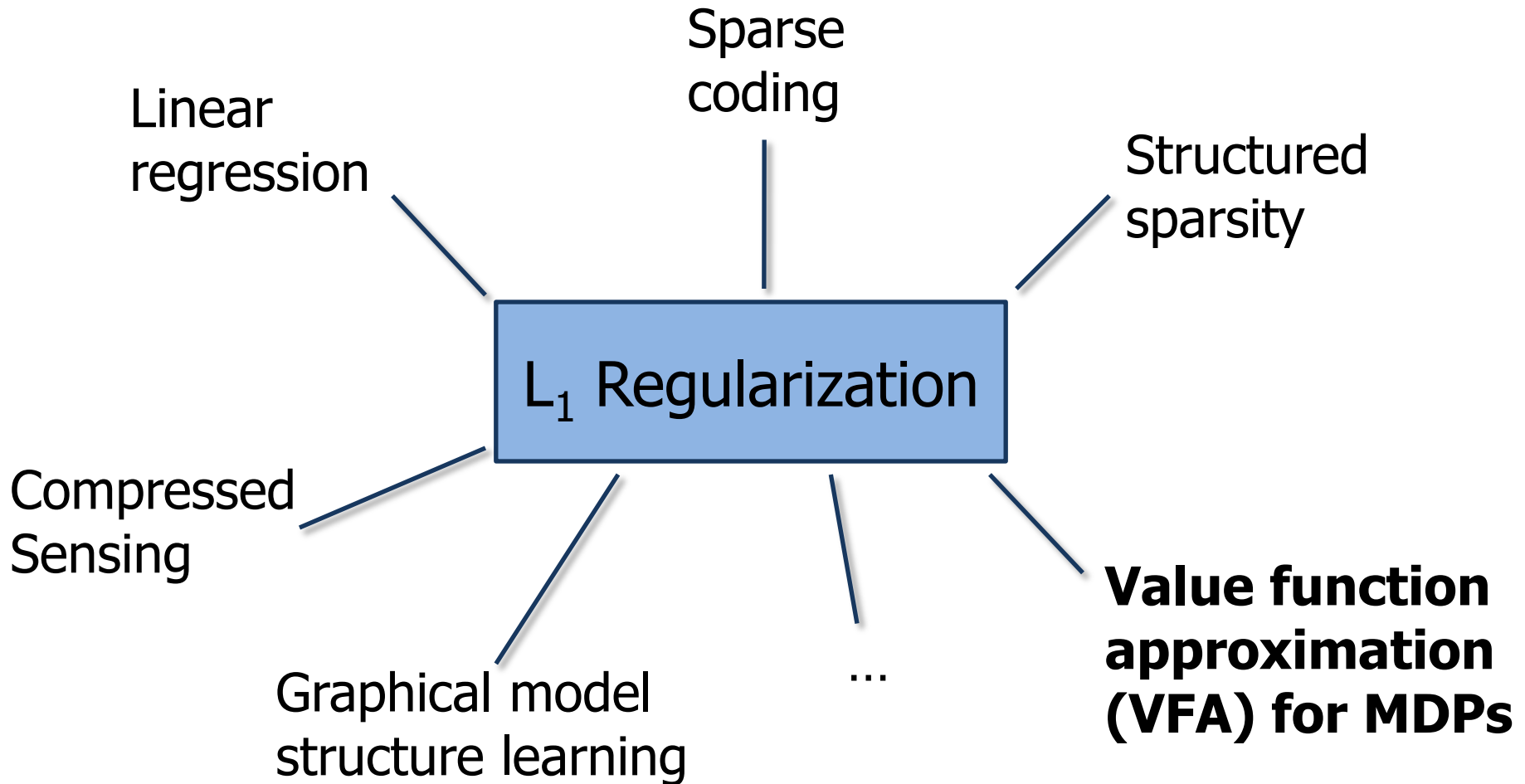
# L$_1$ Regularization

- Form
  - Optimization:  $$w \;=\; \operatorname*{argmin}_{u} \left( L(u) + \beta \|u\|_1 \right)$$

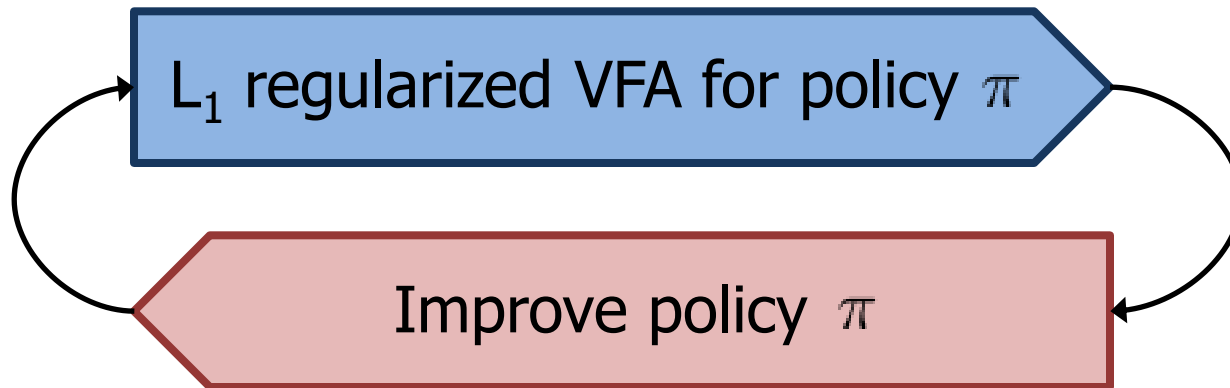  - Fixed point:  $$w \;=\; \operatorname*{argmin}_{u} \left( L(u,w) + \beta \|u\|_1 \right)$$

- Benefits
  - *Sparse* solutions (model interpretability)
  - Learn with large number of irrelevant features
  - Prevents overfitting

# Applications

Sparse
coding

Linear
regression

Structured
sparsity

## $L_1$ Regularization

Compressed
Sensing

Graphical model
structure learning

…

**Value function
approximation
(VFA) for MDPs**

# VFA for MDPs

- Policy iteration



$L_1$ regularized VFA for policy $\pi$

Improve policy $\pi$

- Role of linear complementarity
  - New VFA method
  - Warm starts accelerate computation across iterations

# Agenda

✔ • Big Picture

• Background

• Connection to LCPs

• Two LCP-based Algorithms

• Summary

# Markov Decision Processes

- Finite MDP: $\langle S, A, P, R \rangle$

- Goal: Learn a "good" policy $\pi : S \rightarrow A$

- Value function

$$V^\pi(s) = R(s) + \gamma \sum_{s' \in S} P_{ss'}^{\pi(s)} V^\pi(s')$$

$$V^\pi = R + \gamma P^\pi V^\pi \equiv T^\pi(V^\pi)$$
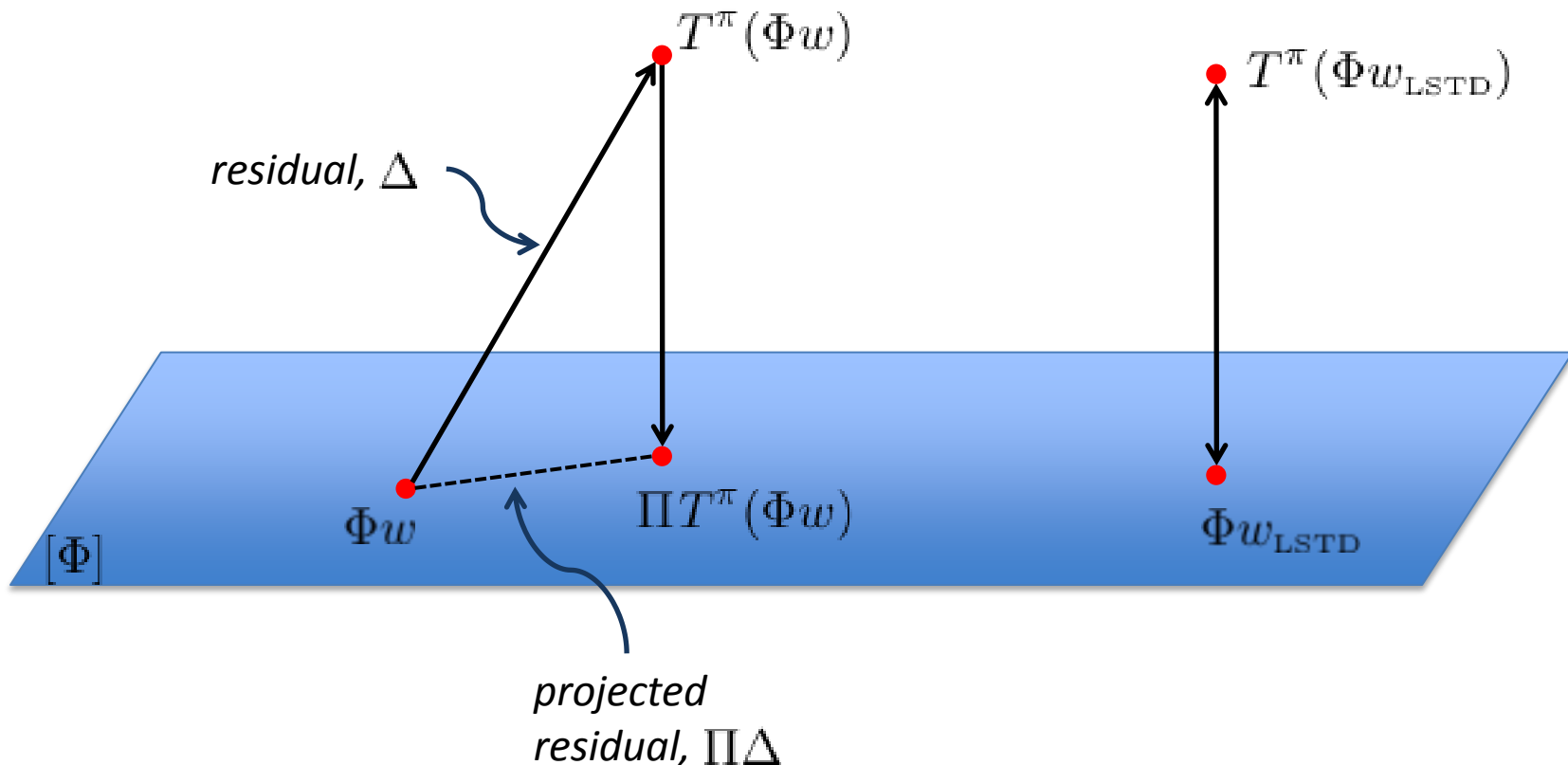
# Linear VFA

- Linear approximation

$$\hat{V}^\pi(s) = \sum_j \phi_j(s)\, w_j$$

$$\hat{V}^\pi = \Phi w$$

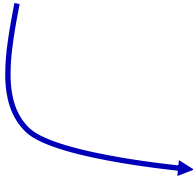- Set $w$ by considering residual $\Delta = T^\pi(\Phi w) - \Phi w$

# LSTD  [Bradtke & Barto, '96]

- Set $w$ to minimize projection of the residual



residual, $\Delta$

$T^\pi(\Phi w)$

$T^\pi(\Phi w_{\mathrm{LSTD}})$

$\Phi w$

$\Pi T^\pi(\Phi w)$

$\Phi w_{\mathrm{LSTD}}$

$[\Phi]$

projected residual, $\Pi\Delta$

# $L_1$TD

- $L_1$TD  =  LSTD + $L_1$ regularization

- $L_1$ fixed point

$$w \;=\; \underset{u}{\operatorname{argmin}} \left( L(u, w) + \beta \|u\|_1 \right)$$

$$\frac{1}{2} \|\Phi u - T^\pi(\Phi w)\|_2^2$$

# L$_1$TD Solution

- No analytical solution

- Coefficients must meet first order conditions

*Derivation parallels*
*LASSO conditions*

$$c = \Phi^T \Delta$$

$$-\beta \leq c_i \leq \beta \qquad \forall i$$
$$c_i = \beta \implies w_i \geq 0$$
$$c_i = -\beta \implies w_i \leq 0$$
$$-\beta < c_i < \beta \implies w_i = 0$$

$\beta \geq 0$ , *L$_1$ reg. parameter*

# LARS-TD   [Kolter & Ng, '09]

- Solves $L_1$TD problem

- Homotopy method inspired by LARS   [Efron et al., '03]
  – Gives solution for *all* values of $\beta$

- Disadvantages w.r.t. policy iteration
  – Each policy starts from scratch
  – Must commit to one $\beta$ to perform policy improvement

# Agenda

✔ • Big Picture

✔ • Background

• Connection to LCPs

• Two LCP-based Algorithms

• Summary

# Linear Complementarity Problems (LCPs)

- LCPs are mathematical programs
  - LP $\subset$ LCP $\subset$ QP

- Given: $\quad q, M$ (square)

  Compute: $\quad d, x$

$$d = Mx + q$$

$$d \geq \mathbf{0}, \ \ x \geq \mathbf{0}$$

$$d_i x_i = 0 \quad \forall i$$

# L$_1$TD as an LCP

- Recall  $c \ =\ \Phi^T \Delta \ =\ \underbrace{\Phi^T R}_{b} - \underbrace{\Phi^T(\Phi - \gamma P^\pi \Phi)w}_{A}$

- LCP inputs:  $q \ =\ \left(\beta + \begin{bmatrix} b \\ -b \end{bmatrix}\right)$

$$M \ =\ \begin{bmatrix} A & -A \\ -A & A \end{bmatrix}$$

- LCP solution achieves L$_1$TD first order conditions
  - Denote this as  $w \leftarrow \text{LC-TD}(A, b, \beta)$

# LCP Solvers

- Similarities to LP solvers
  - Worst case exponential complexity

- Some solvers can be initialized with a *warm start*

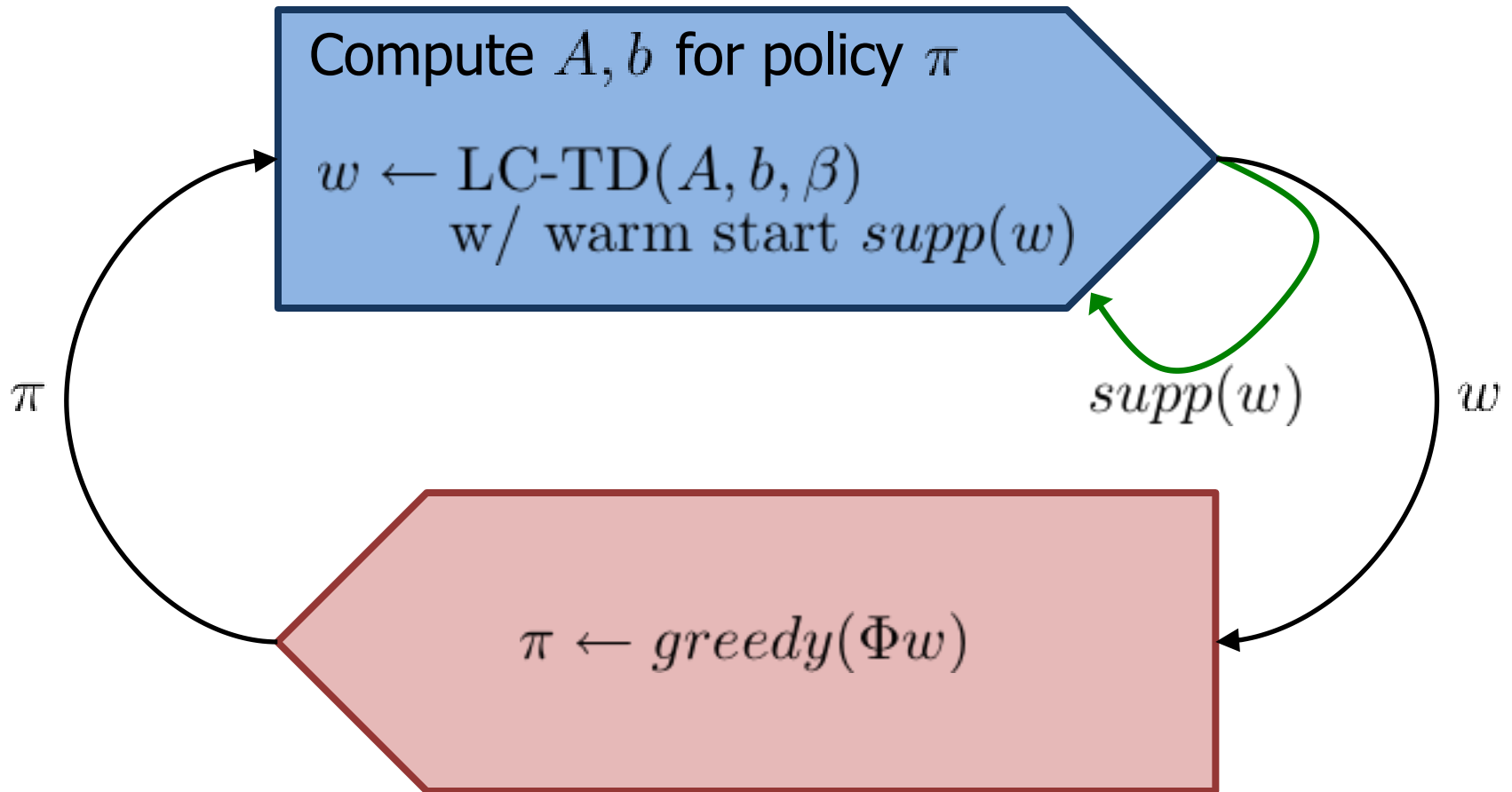- Properties of $q, M$ dictate solution existence, uniqueness, and computability

# Agenda

✔ • Big Picture

✔ • Background

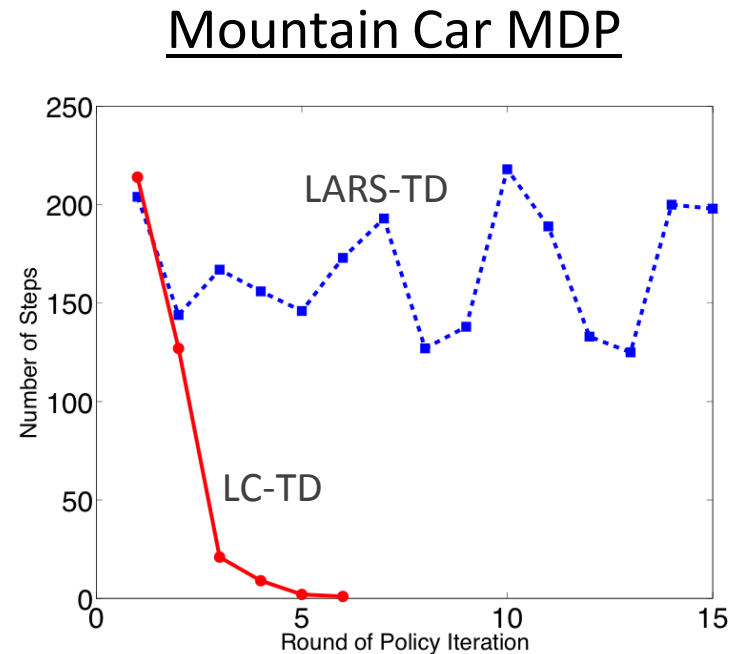✔ • Connection to LCPs

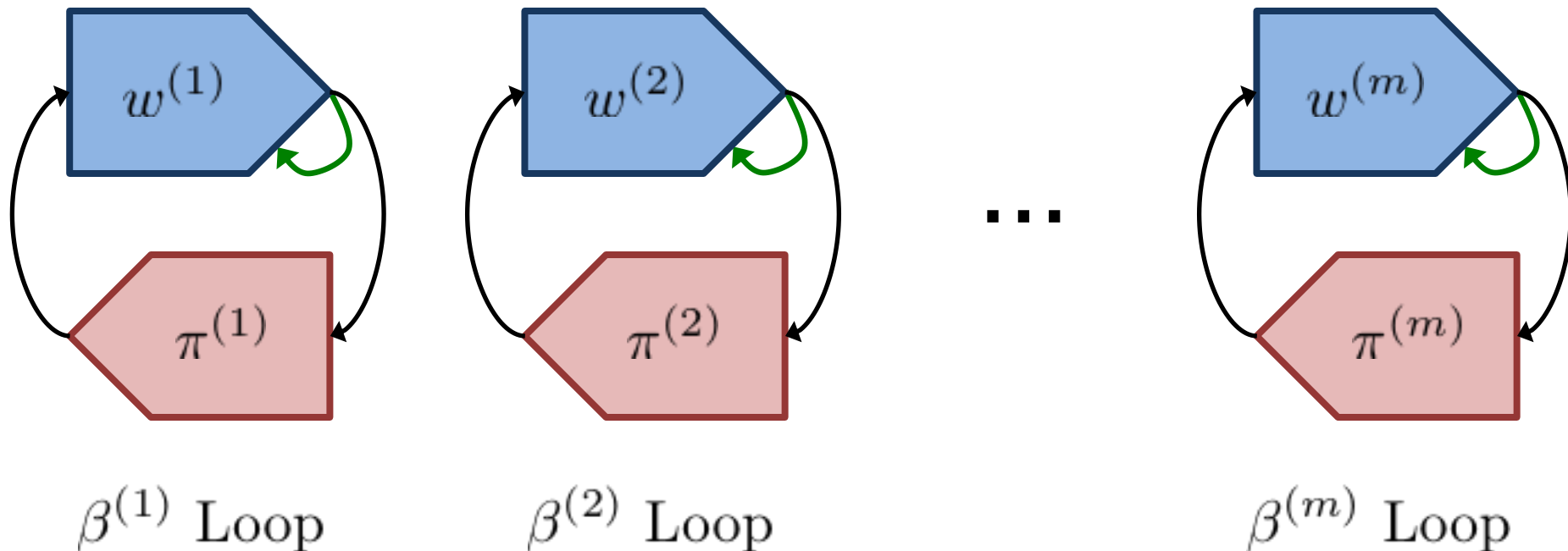• Two LCP-based Algorithms

• Summary

# Algorithm 1: LC-TD w/ PI



Compute $A, b$ for policy $\pi$

$w \leftarrow \text{LC-TD}(A, b, \beta)$
w/ warm start $supp(w)$

$supp(w)$

$w$

$\pi$

$\pi \leftarrow greedy(\Phi w)$

# Experiments

- Compare average number of algorithm steps for LC-TD and LARS-TD

| Domain | LARS-TD, PI | LC-TD, PI |
|---|---|---|
| Mountain Car | 214 ± 23 | **116 ± 22** |
| Chain | 73 ± 13 | 77 ± 11 |
| Pendulum | 153 ± 25 | **48 ± 23** |

Mountain Car MDP

# Choosing $\beta$

- Option 1: Run $m$ separate trials



$\beta^{(1)}$ Loop $\quad\quad$ $\beta^{(2)}$ Loop $\quad\quad\quad$ $\beta^{(m)}$ Loop
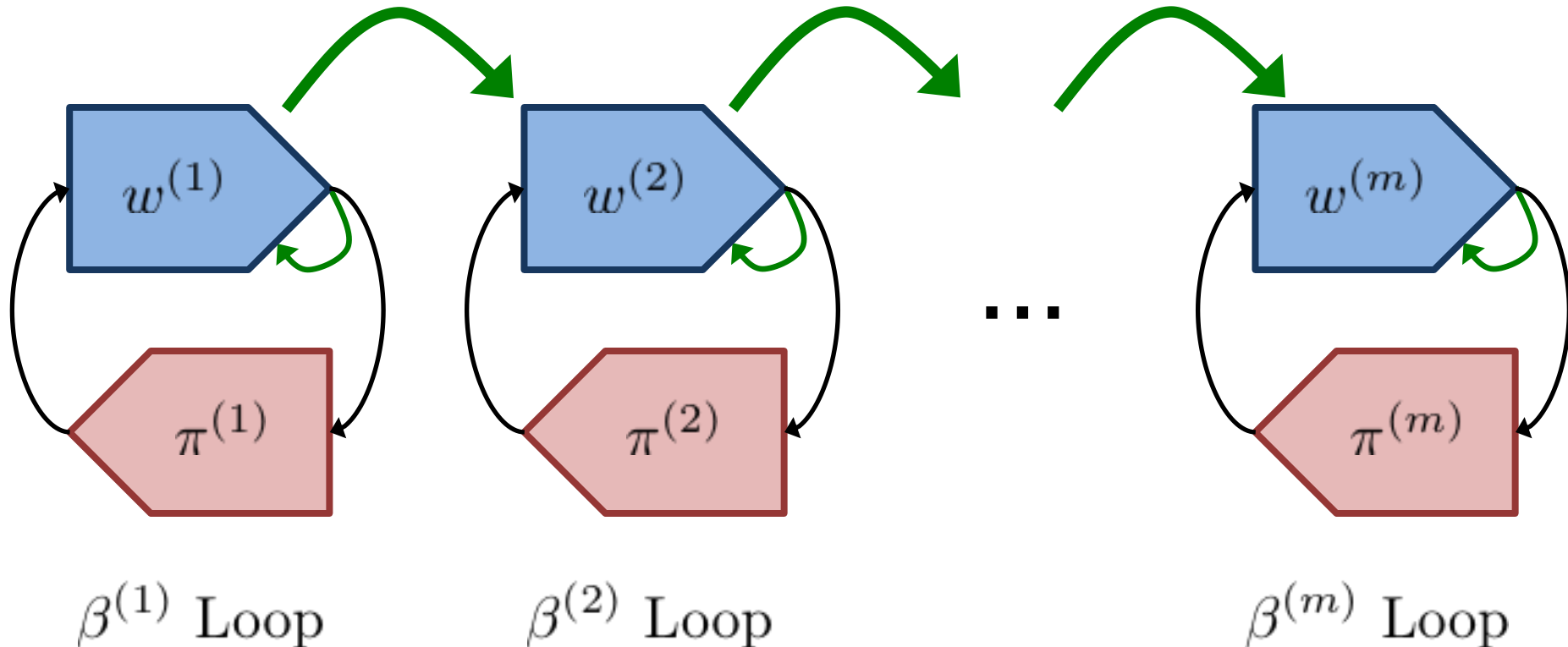
# Algorithm 2: LC-MPI

Let $\beta^{(1)} > \beta^{(2)} > \ldots > \beta^{(m)}$

Initialize:

$$w^{(2)} \leftarrow w^{(1)}$$
$$\pi^{(2)} \leftarrow \pi^{(1)}$$



$w^{(1)}$    $w^{(2)}$    $\ldots$    $w^{(m)}$

$\pi^{(1)}$    $\pi^{(2)}$    $\pi^{(m)}$

$\beta^{(1)}$ Loop    $\beta^{(2)}$ Loop    $\beta^{(m)}$ Loop

# Experiments

- Compare avg. number of algorithm steps
  - LARS-TD and LC-TD 11 separate times for different $\beta$'s
  - LC-MPI once for same 11 $\beta$'s

| Domain | LARS-TD, PI | LC-TD, PI | LC-MPI |
|--------|-------------|-----------|--------|
| Mountain Car | 214 ± 23 | 116 ± 22 | **21 ± 5** |
| Chain | 73 ± 13 | 77 ± 11 | **24 ± 11** |
| Pendulum | 153 ± 25 | 48 ± 23 | **33 ± 18** |

# Agenda

✔ • Big Picture

✔ • Background

✔ • Connection to LCPs

✔ • Two LCP-based Algorithms

• Summary

# Summary

- $L_1$ fixed point is equivalent to an LCP

- LCP formulation useful when solving multiple, related problems
  - Exploit warm starts
  - Reduce computation during policy iteration
  - LC-MPI lends itself to cross-validation

# Thanks