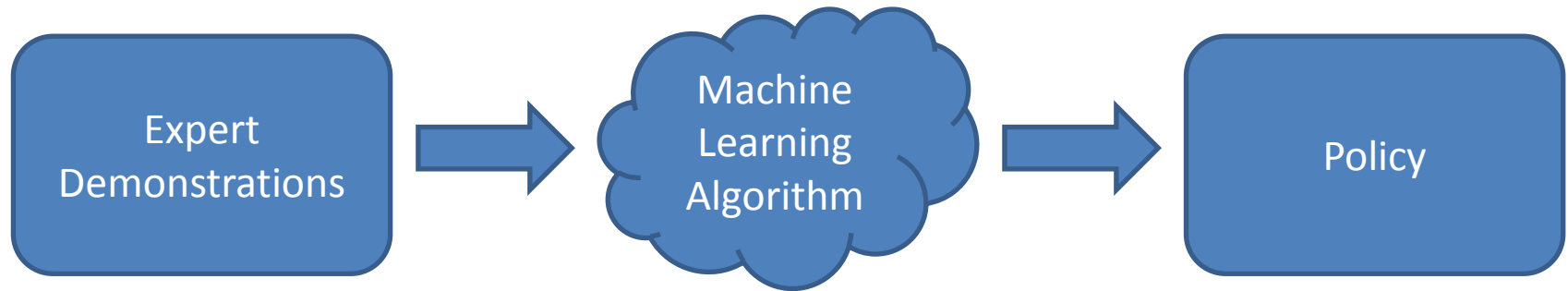# Reduction of Imitation Learning to No-Regret Online Learning

Stephane Ross

Joint work with Drew Bagnell & Geoff Gordon

**Carnegie Mellon**
**THE ROBOTICS INSTITUTE**

# Imitation Learning



Expert Demonstrations → Machine Learning Algorithm → Policy

# Imitation Learning

- Many successes:
  - Legged locomotion [Ratliff 06]
  - Outdoor navigation [Silver 08]
  - Helicopter flight [Abbeel 07]
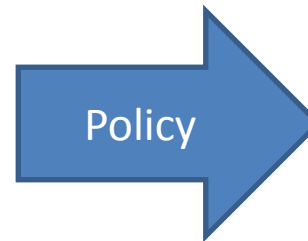  - Car driving [Pomerleau 89]
  - etc...

# Example Scenario
## Learning to drive from demonstrations

Input:                                    Output:
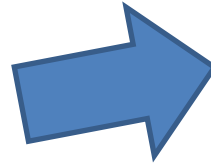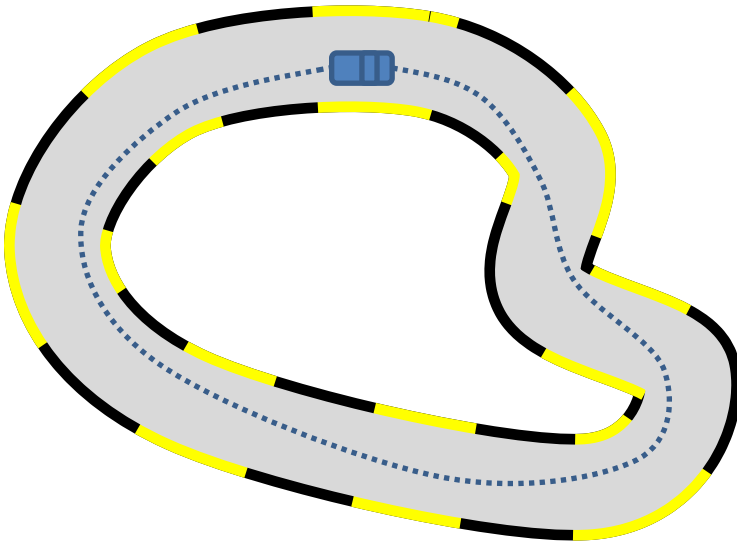


Camera Image

Policy

Steering in [-1,1]

Hard left turn          Hard right turn

# Supervised Training Procedure

Expert Trajectories

Dataset



Learned Policy: $\hat{\pi}_{\mathbf{sup}} = \mathbf{arg\,min}_{\pi \in \Pi} \mathop{\mathbf{E}}_{\mathbf{s} \sim \mathbf{D}(\pi^*)} [\ell(\pi, \mathbf{s}, \pi^*(\mathbf{s}))]$

# Poor Performance in Practice

# # Mistakes Grows Quadratically in T!

[Ross 2010]

$$\mathbf{J}(\hat{\pi}_{\mathbf{sup}}) \leq \mathbf{T}^2 \varepsilon$$
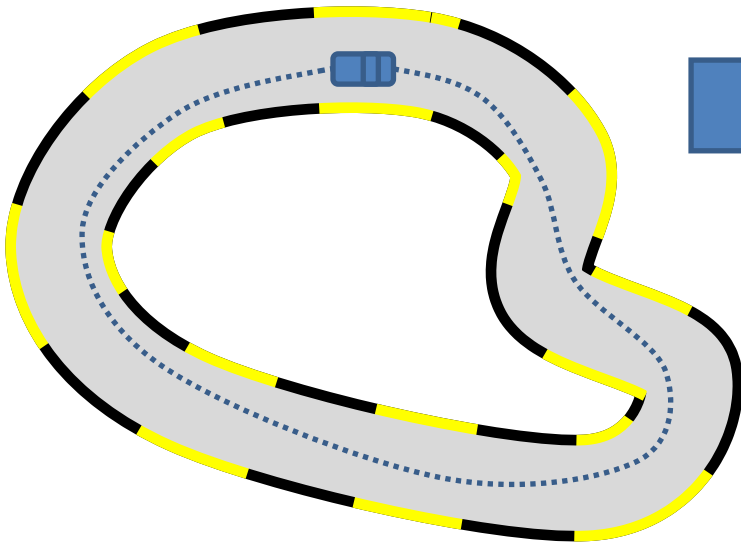
Exp. # of mistakes over T steps

# time steps

Avg. loss on $D(\pi*)$
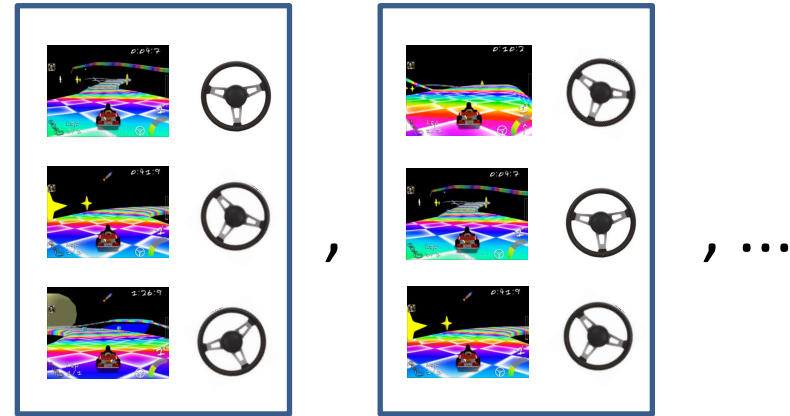
**Reason: Doesn't learn how to recover from errors!**

# Reduction-Based Approach & Analysis

**Hard Learning Problem**

**Easier Related Problem(s)**



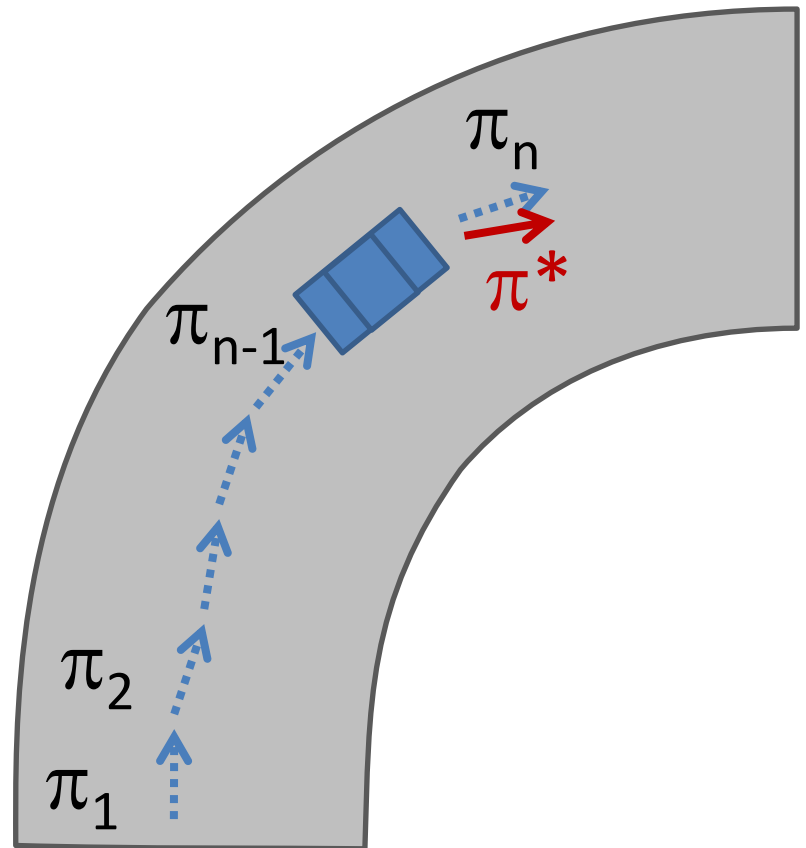Performance: f(ε)

Performance: ε

**Example:** Cost-sensitive Multiclass classification to Binary classification [Beygelzimer 2005]

# Previous Work: Forward Training

[Ross 2010]

- Sequentially learn one policy/step

- # mistakes grows linearly:
  - $J(\pi_{1:T}) \leq T\varepsilon$

- **Impractical if T large**

$\pi_n$

$\pi^*$

$\pi_{n-1}$
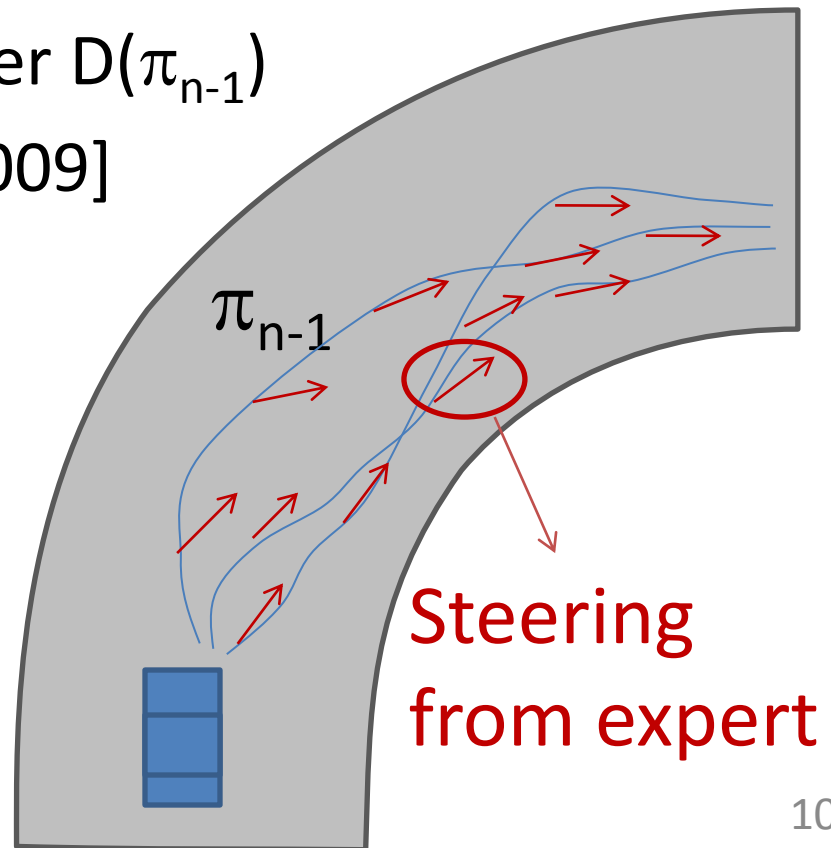
$\pi_2$

$\pi_1$

# Previous Work: SMILe

- Learn stochastic policy, changing policy slowly
  - $\pi_n = \pi_{n-1} + \alpha_n(\pi'_n - \pi^*)$
  - $\pi'_n$ trained to mimic $\pi^*$ under $D(\pi_{n-1})$
  - Similar to SEARN [Daume 2009]
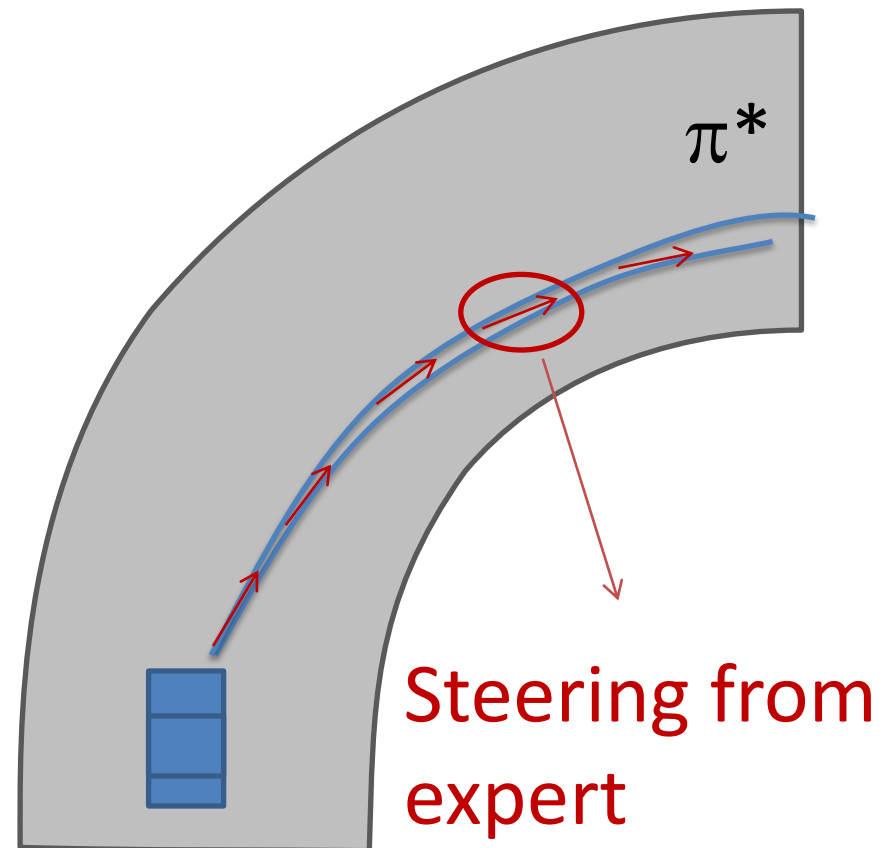
- Near-linear bound:
  - $J(\pi) \leq O(T\log(T)\varepsilon + 1)$
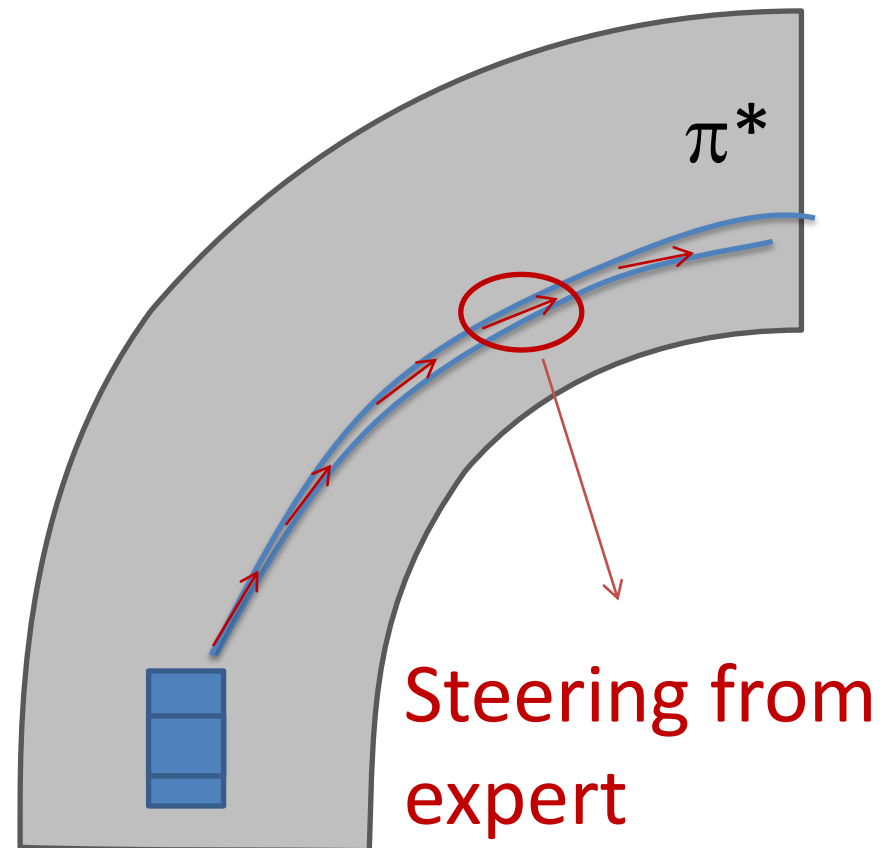
- **Stochasticity undesirable**

$\pi_{n-1}$

Steering
from expert

# DAgger: Dataset Aggregation

- Collect trajectories with expert $\pi^*$

$\pi^*$

Steering from expert

# DAgger: Dataset Aggregation

- Collect trajectories with expert $\pi^*$
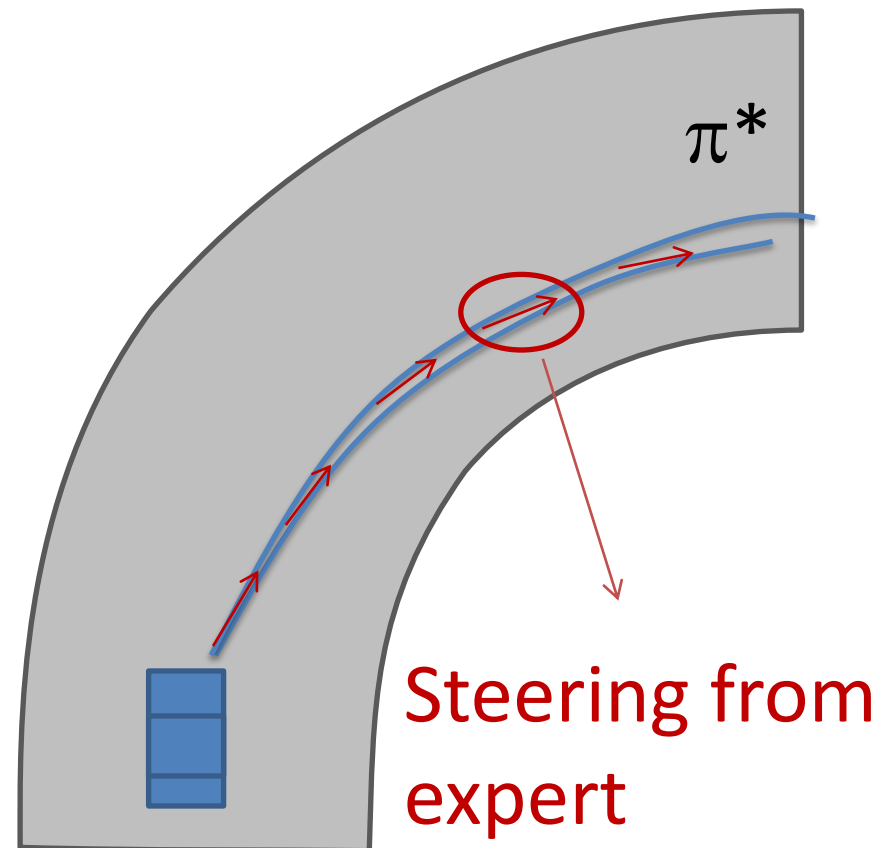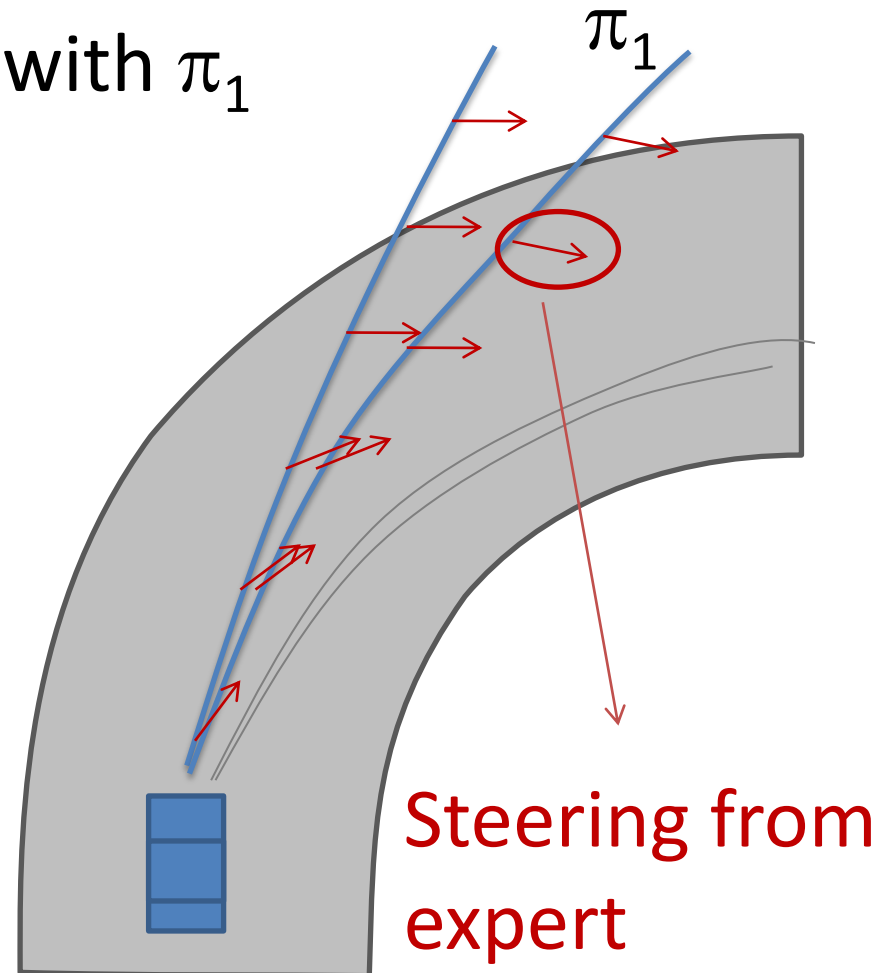
- Dataset $D_0 = \{(s, \pi^*(s))\}$

$\pi^*$

Steering from expert

# DAgger: Dataset Aggregation

- Collect trajectories with expert $\pi^*$

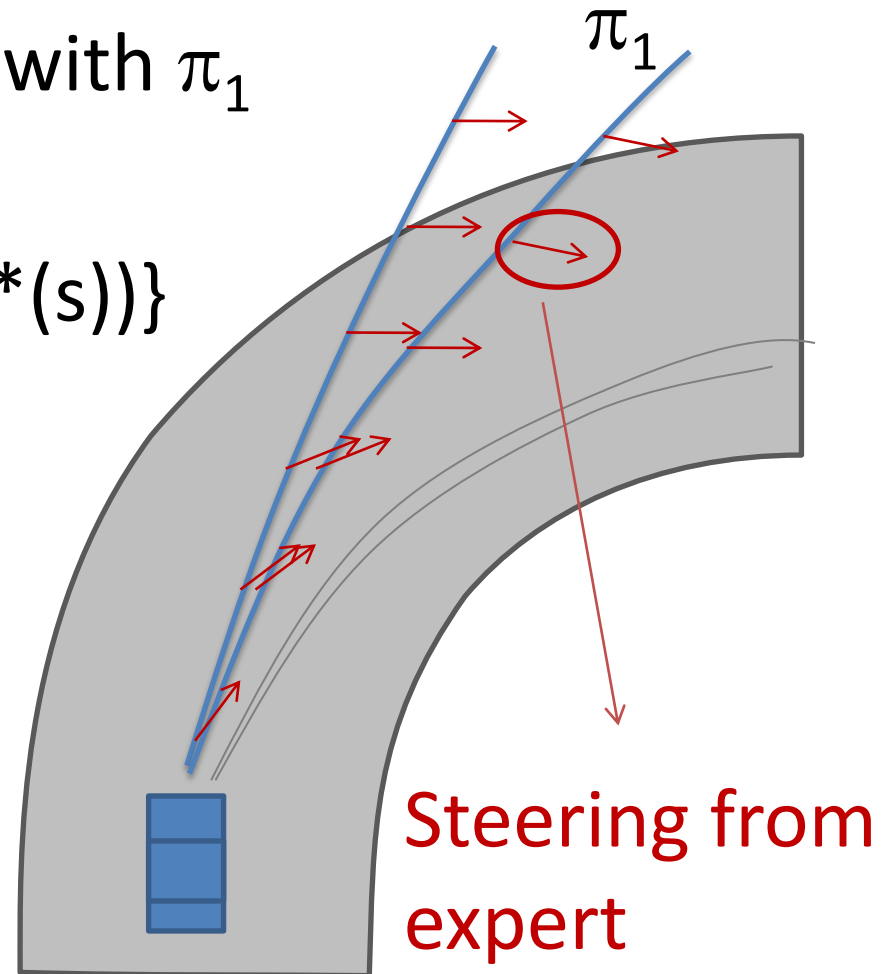- Dataset $D_0 = \{(s, \pi^*(s))\}$

- Train $\pi_1$ on $D_0$

$\pi^*$

Steering from expert

# DAgger: Dataset Aggregation

- Collect new trajectories with $\pi_1$

$\pi_1$

Steering from expert

# DAgger: Dataset Aggregation

- Collect new trajectories with $\pi_1$

- New Dataset $D_1' = \{(s, \pi^*(s))\}$
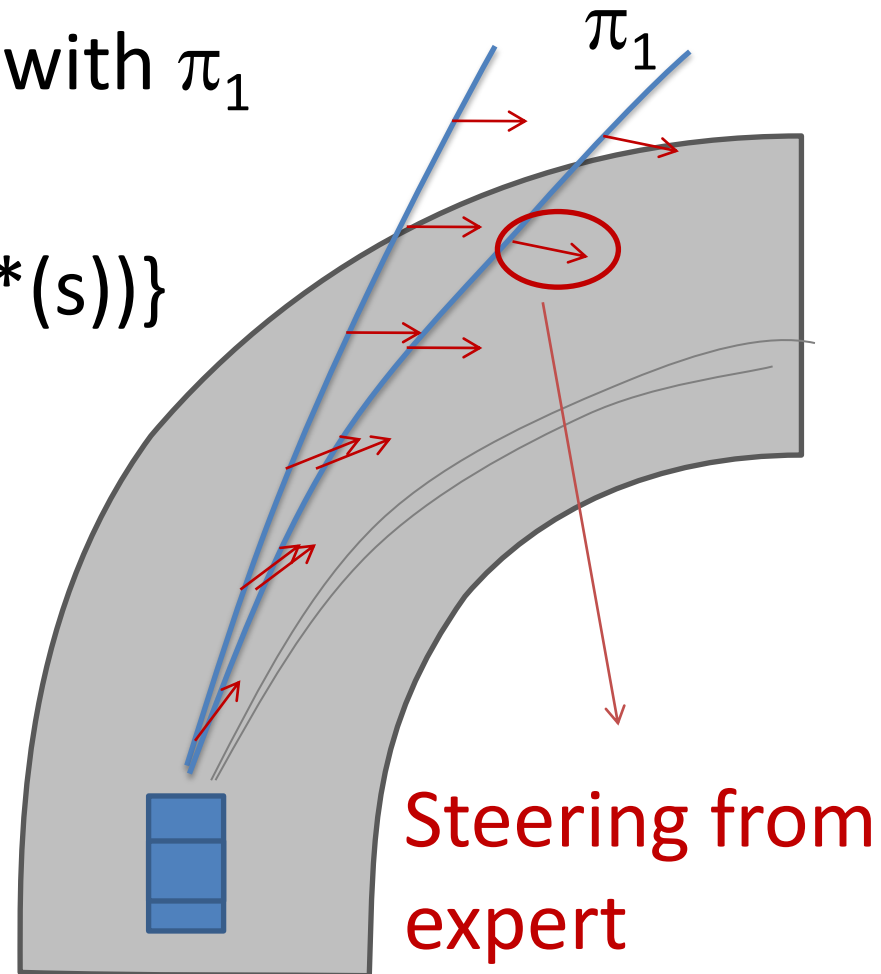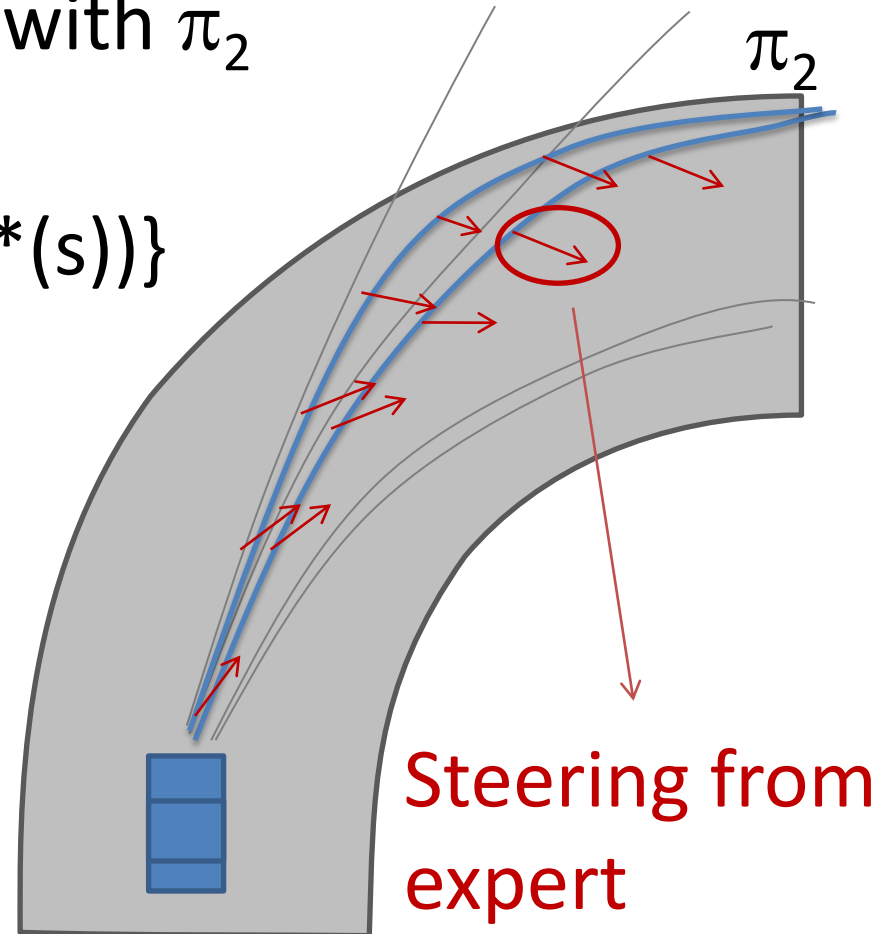
$\pi_1$

Steering from expert

# DAgger: Dataset Aggregation

- Collect new trajectories with $\pi_1$

- New Dataset $D_1' = \{(s, \pi^*(s))\}$

- Aggregate Datasets:
  $D_1 = D_0 \cup D_1'$

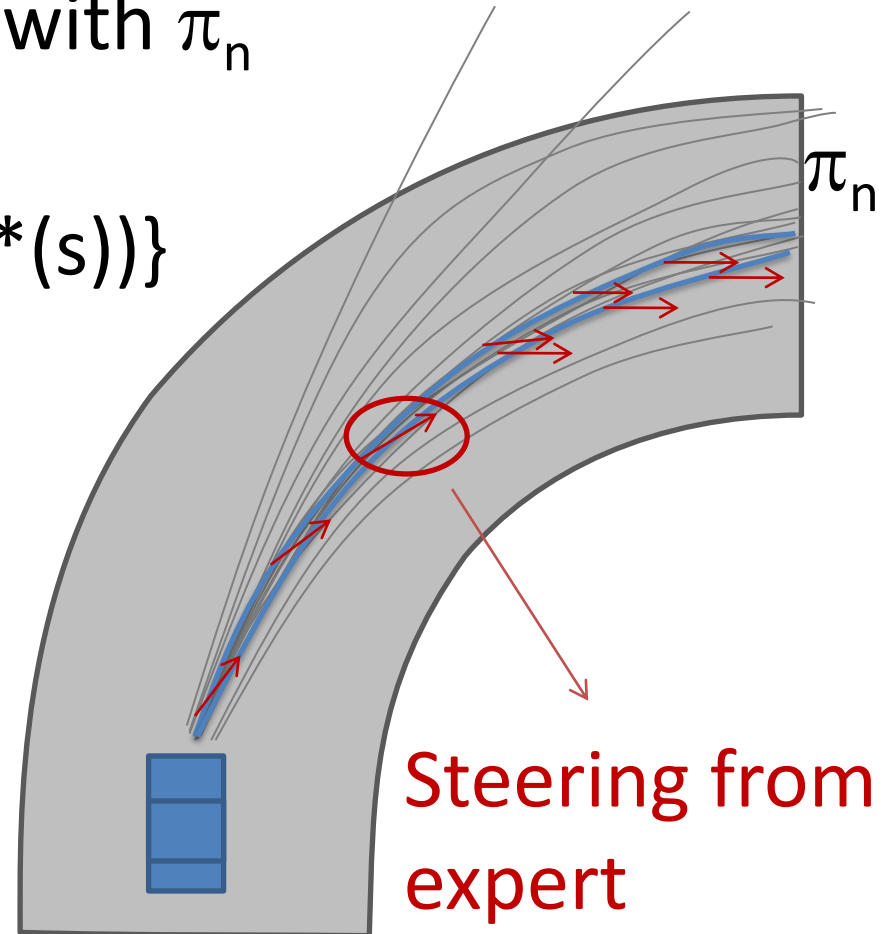

$\pi_1$

Steering from expert

# DAgger: Dataset Aggregation

- Collect new trajectories with $\pi_1$

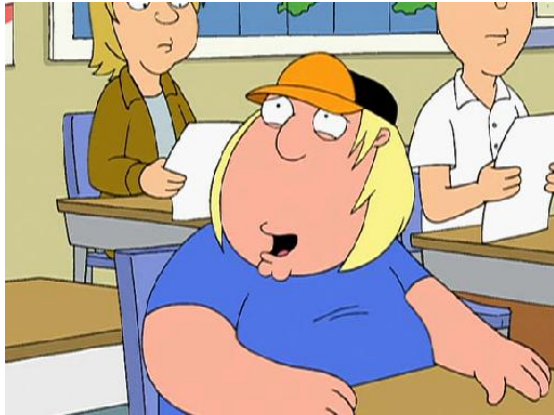- New Dataset $D_1' = \{(s, \pi^*(s))\}$

- Aggregate Datasets:
  $D_1 = D_0 \cup D_1'$

- Train $\pi_2$ on $D_1$

$\pi_1$

Steering from expert

# DAgger: Dataset Aggregation

- Collect new trajectories with $\pi_2$

- New Dataset $D_2' = \{(s, \pi^*(s))\}$

- Aggregate Datasets:
  $D_2 = D_1 \cup D_2'$

- Train $\pi_3$ on $D_2$

$\pi_2$

Steering from expert

# DAgger: Dataset Aggregation

- Collect new trajectories with $\pi_n$

- New Dataset $D_n' = \{(s, \pi^*(s))\}$

- Aggregate Datasets:
  $D_n = D_{n-1} \cup D_n'$
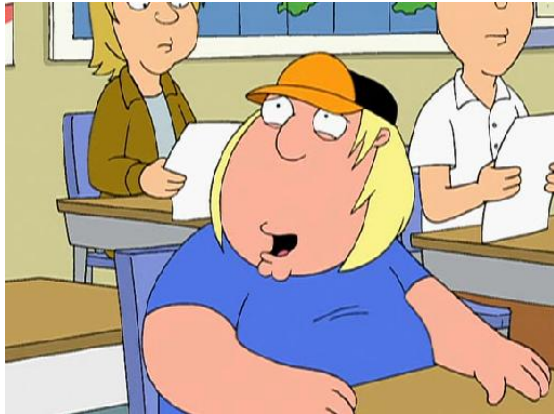
- Train $\pi_{n+1}$ on $D_n$

$\pi_n$

Steering from expert
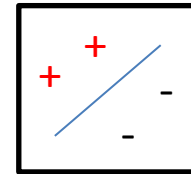
# Online Learning

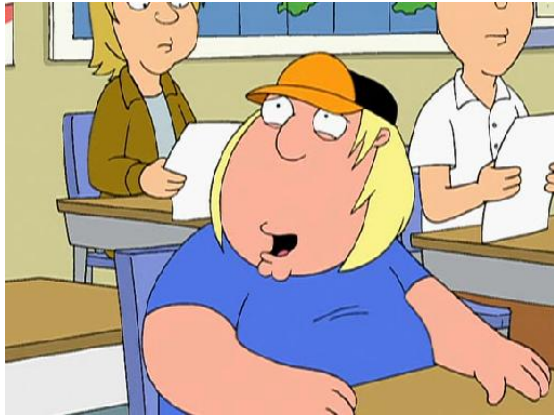**Learner**

**Adversary**

# Online Learning
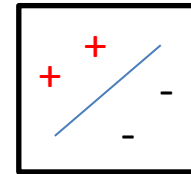
**Learner**

**Adversary**

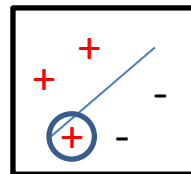Current Hypothesis $h_n$

# Online Learning
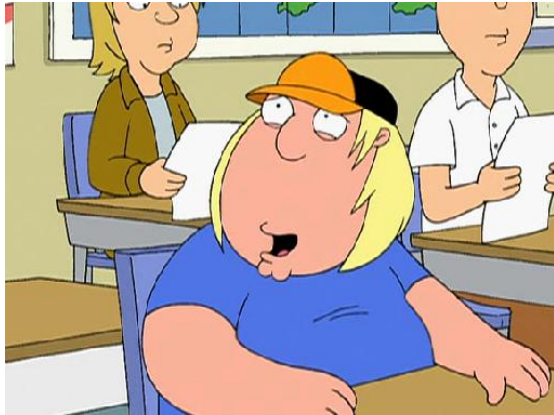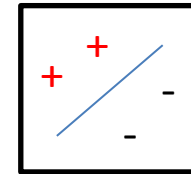
**Learner**

**Adversary**

Current Hypothesis $h_n$
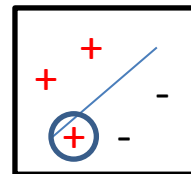
Pick Loss $L_n$
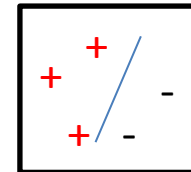
# Online Learning
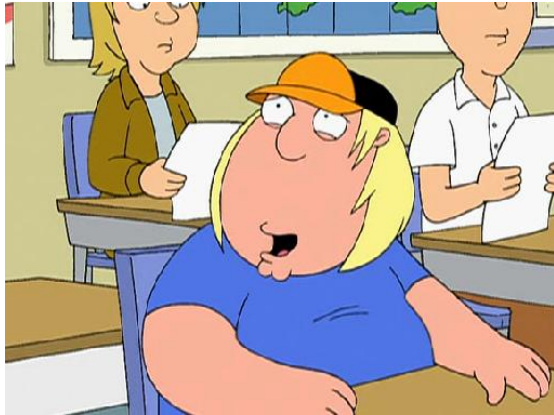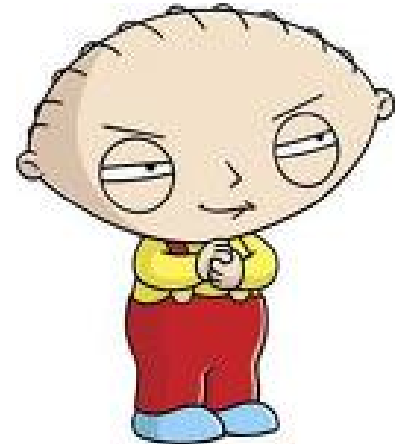
**Learner**

**Adversary**

Current Hypothesis $h_n$

Pick Loss $L_n$

Next Hypothesis $h_{n+1}$

# Online Learning

**Learner**

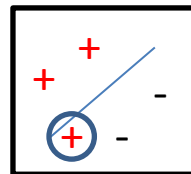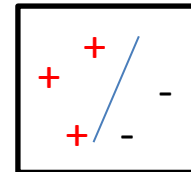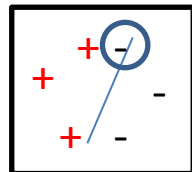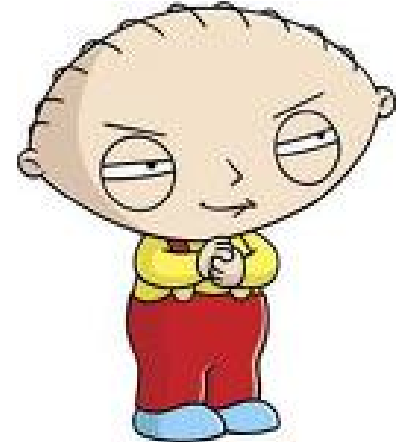**Adversary**

# Online Learning



**Learner**

**Adversary**

$\vdots$

Current Hypothesis $h_n$

Pick Loss $L_n$

Next Hypothesis $h_{n+1}$

Pick Loss $L_{n+1}$

Avg. Regret: $\gamma_n = \dfrac{1}{n}\left[\sum_{i=1}^{n} L_i(h_i) - \min_{h \in H} \sum_{i=1}^{n} L_i(h)\right]$

# DAgger as Online Learning
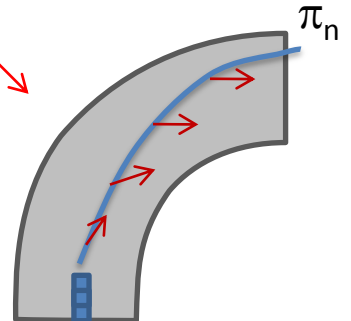


**Learner**

**Adversary**

⋮

Current Policy $\pi_n$

Pick Loss $L_n$

Next Policy $\pi_{n+1}$

Pick Loss $L_{n+1}$

$$\mathbf{L_n}(\pi) = \underset{\mathbf{s \sim D(\pi_n)}}{\mathbf{E}} [\ell(\pi, \mathbf{s}, \pi^*(\mathbf{s}))]$$

26

# DAgger as Online Learning



$$\pi_{\mathbf{n+1}} = \underset{\pi \in \Pi}{\arg\min} \sum_{\mathbf{i}=1}^{\mathbf{n}} \mathbf{L_i}(\pi)$$

$$\mathbf{L_n}(\pi) = \underset{\mathbf{s} \sim \mathbf{D}(\pi_{\mathbf{n}})}{\mathbf{E}} [\ell(\pi, \mathbf{s}, \pi^*(\mathbf{s}))]$$

# DAgger as Online Learning



**Learner**
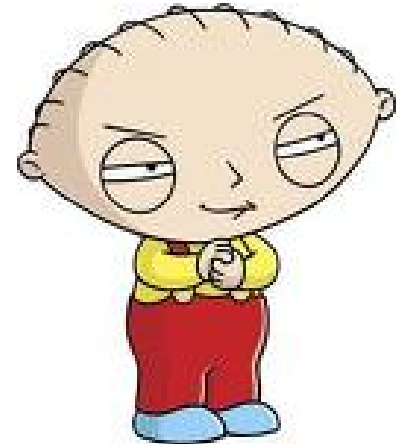
$\vdots$

**Adversary**

Current Policy $\pi_n$
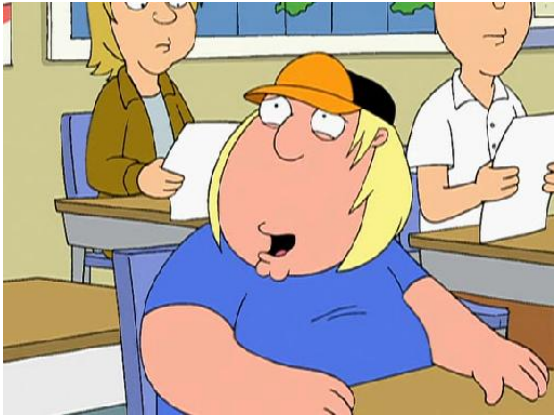
Pick Loss $L_n$

Next Policy $\pi_{n+1}$

$$\pi_{n+1} = \underset{\pi \in \Pi}{\arg\min} \sum_{i=1}^{n} L_i(\pi)$$

Pick Loss $L_{n+1}$

$\pi_n$

Follow-The-Leader (FTL)

$$L_n(\pi) = \underset{s \sim D(\pi_n)}{E}[\ell(\pi, s, \pi^*(s))]$$

28

# Theoretical Guarantees of DAgger

- Best policy $\pi$ in sequence $\pi_{1:N}$ guarantees:

$$\mathbf{J}(\pi) \leq \mathbf{T}(\varepsilon_{\mathbf{N}} + \gamma_{\mathbf{N}}) + \mathbf{O}(\mathbf{T}/\mathbf{N})$$

Avg. Loss on Aggregate Dataset

Avg. Regret of $\pi_{1:N}$

Iterations of DAgger

# Theoretical Guarantees of DAgger

- Best policy $\pi$ in sequence $\pi_{1:N}$ guarantees:

$$\mathbf{J}(\pi) \leq \mathbf{T}(\varepsilon_{\mathbf{N}} + \gamma_{\mathbf{N}}) + \mathbf{O}(\mathbf{T}/\mathbf{N})$$

Avg. Loss on Aggregate Dataset

Avg. Regret of $\pi_{1:N}$

Iterations of DAgger

- For strongly convex loss, **N = O(TlogT) iterations**:

$$\mathbf{J}(\pi) \leq \mathbf{T}\varepsilon_{\mathbf{N}} + \mathbf{O}(1)$$

# Theoretical Guarantees of DAgger

- Best policy $\pi$ in sequence $\pi_{1:N}$ guarantees:

$$\mathbf{J}(\pi) \leq \mathbf{T}(\varepsilon_\mathbf{N} + \gamma_\mathbf{N}) + \mathbf{O}(\mathbf{T}/\mathbf{N})$$

Avg. Loss on Aggregate Dataset

Avg. Regret of $\pi_{1:N}$

Iterations of DAgger

- For strongly convex loss, **N = O(TlogT) iterations**:

$$\mathbf{J}(\pi) \leq \mathbf{T}\varepsilon_\mathbf{N} + \mathbf{O}(1)$$

- Any No-Regret algorithm has same guarantees

# Theoretical Guarantees of DAgger

- If sample **m trajectories** at each iteration, w.p. 1-$\delta$:

$$\mathbf{J}(\pi) \leq \mathbf{T}(\hat{\varepsilon}_{\mathbf{N}} + \gamma_{\mathbf{N}}) + \mathbf{O}(\mathbf{T}\sqrt{\mathbf{log}(1/\delta)}/\sqrt{\mathbf{Nm}})$$

Empirical Avg. Loss on Aggregate Dataset

Avg. Regret of $\pi_{1:N}$

# Theoretical Guarantees of DAgger

- If sample **m trajectories** at each iteration, w.p. 1-$\delta$:

$$\mathbf{J}(\pi) \leq \mathbf{T}(\hat{\varepsilon}_{\mathbf{N}} + \gamma_{\mathbf{N}}) + \mathbf{O}(\mathbf{T}\sqrt{\log(1/\delta)}/\sqrt{\mathbf{Nm}})$$

Empirical Avg. Loss on Aggregate Dataset

Avg. Regret of $\pi_{1:N}$

- For strongly convex loss, **N = O(T²log(1/$\delta$)) , m=1**, w.p. 1-$\delta$:

$$\mathbf{J}(\pi) \leq \mathbf{T}\hat{\varepsilon}_{\mathbf{N}} + \mathbf{O}(1)$$
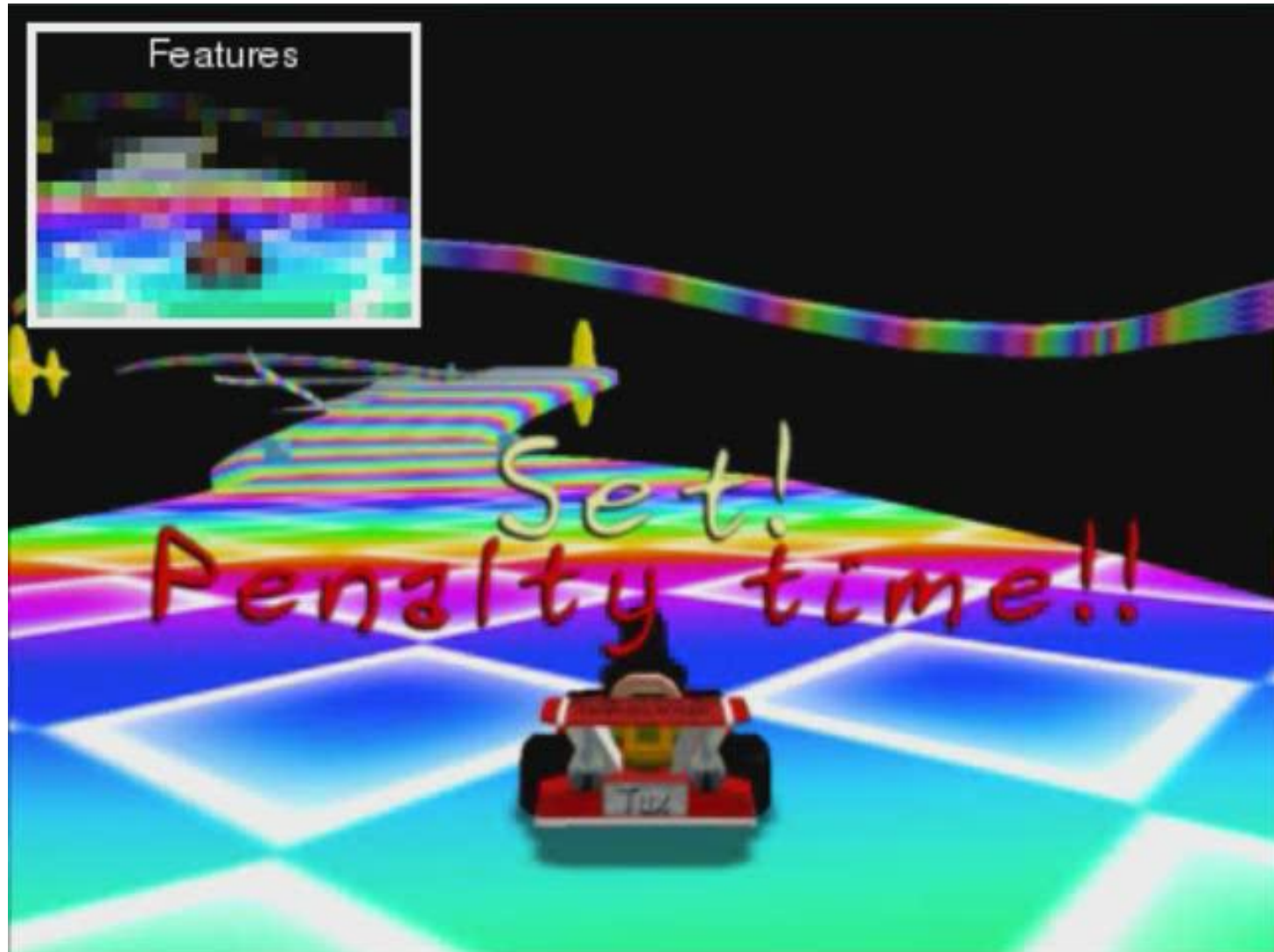
# Experiments: 3D Racing Game

Input:

Output:



Resized to 25x19
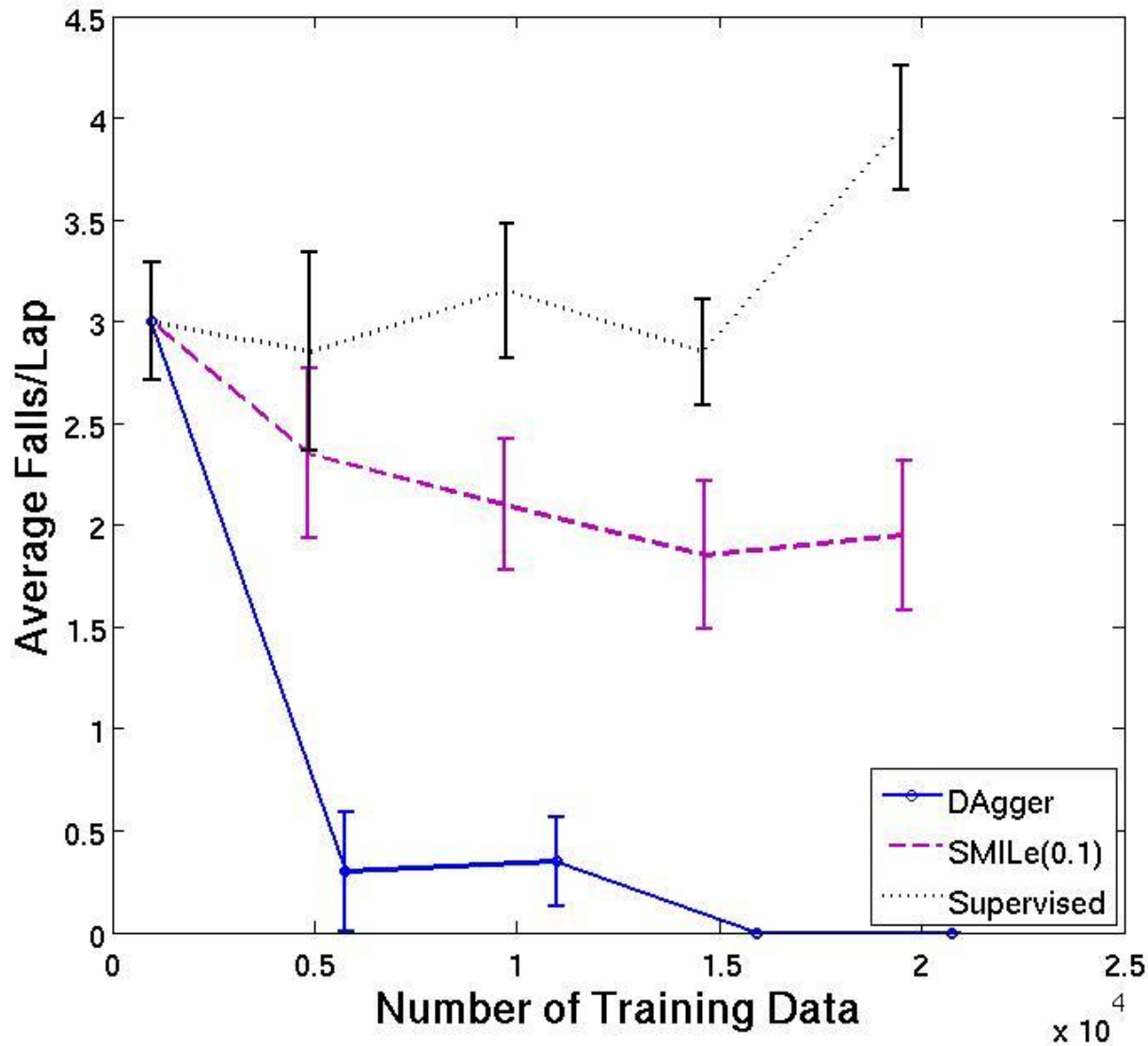pixels (1425 features)

Steering in [-1,1]

# DAgger Test-Time Execution

# Average Falls/Lap
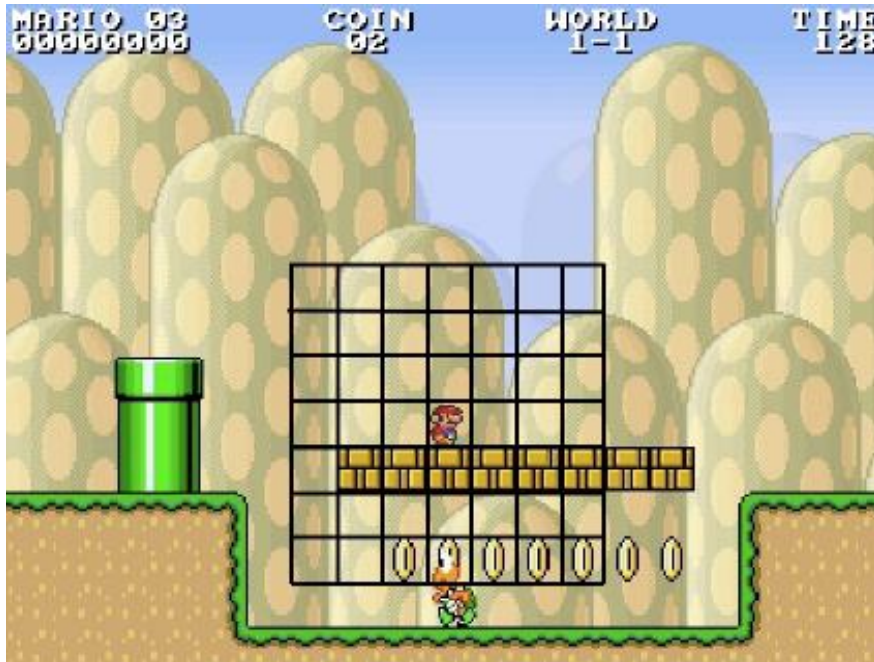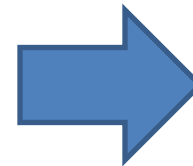
# Experiments: Super Mario Bros
## From Mario AI competition 2009

Input:

Output:



Jump in {0,1}
Right in {0,1}
Left in {0,1}
Speed in {0,1}

Extracted 27K+ binary features from last 4 observations
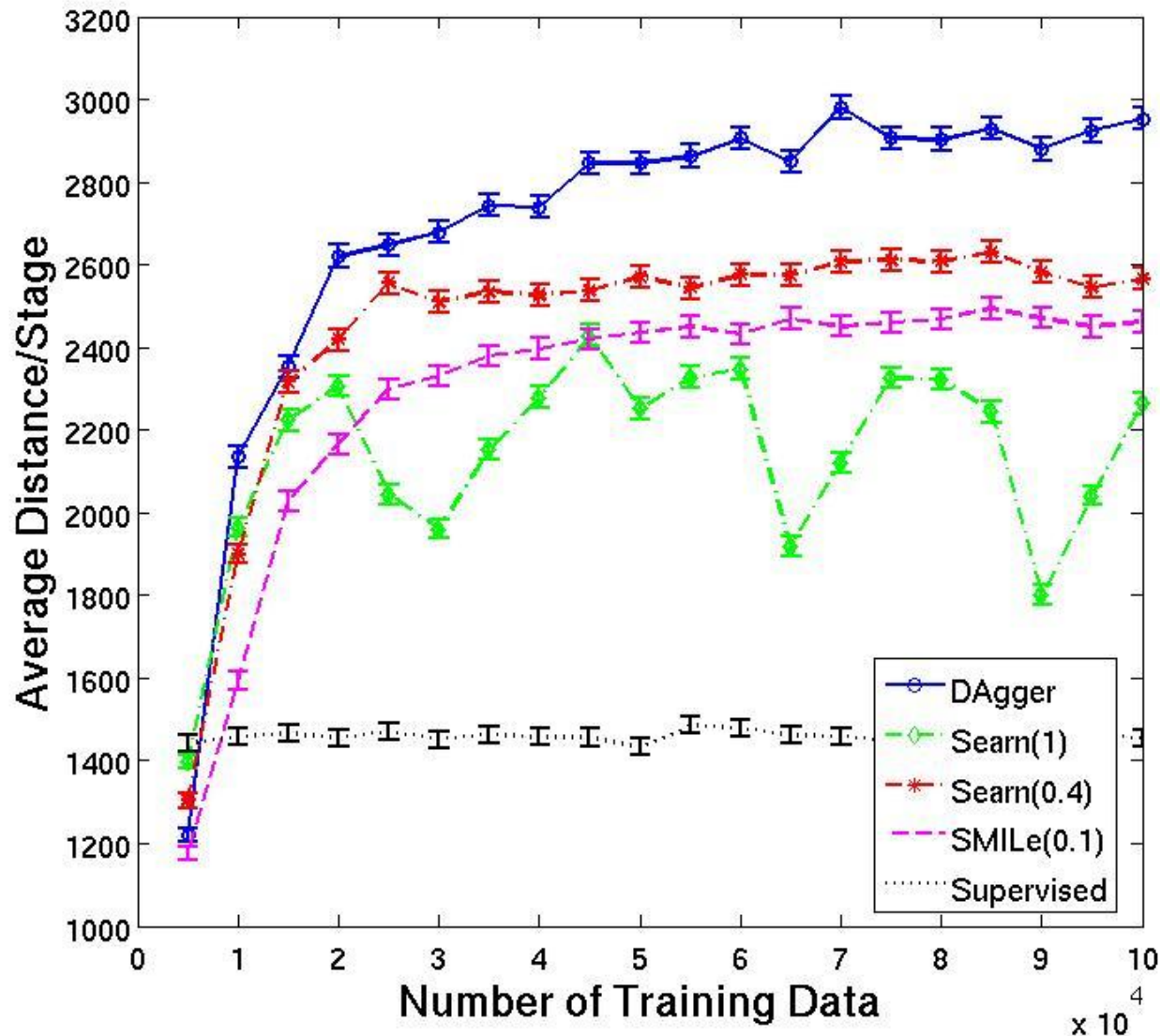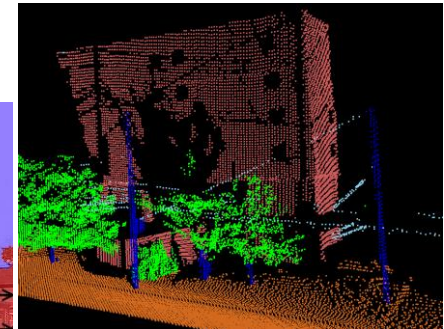(14 binary features for every cell)

# Test-Time Execution

# Average Distance/Stage

# Conclusion

- Take-Home Message
  - Simple iterative procedures can yield much better performance.

- Can also be applied for **Structured Prediction**:
  - NLP (e.g. Handwriting Recognition)
  - Computer Vision [Ross & al., CVPR 2011]

- Future Work:
  - Combining with other Imitation Learning techniques [Ratliff 06]
  - Potential extensions to Reinforcement Learning?
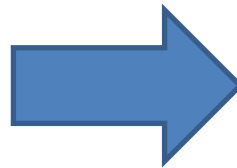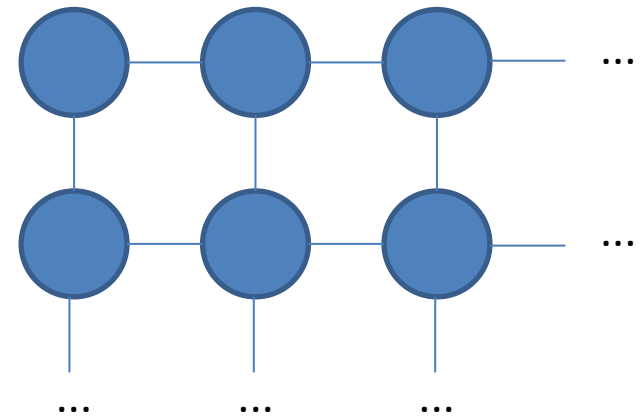
# Questions

# Structured Prediction
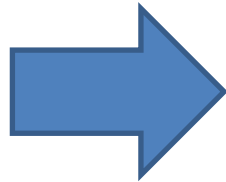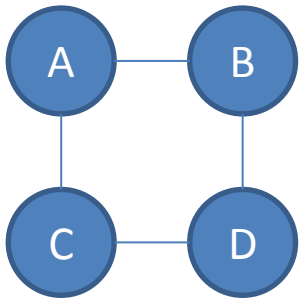
- Example: Scene Labeling

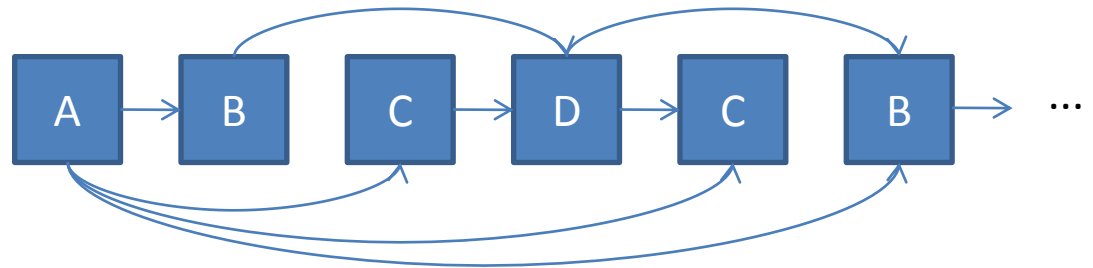Image



Graph Structure over Labels

# Structured Prediction

- Sequentially label each node using neighboring predictions
  - e.g. In Breath-First-Search Order (Forward & Backward passes)

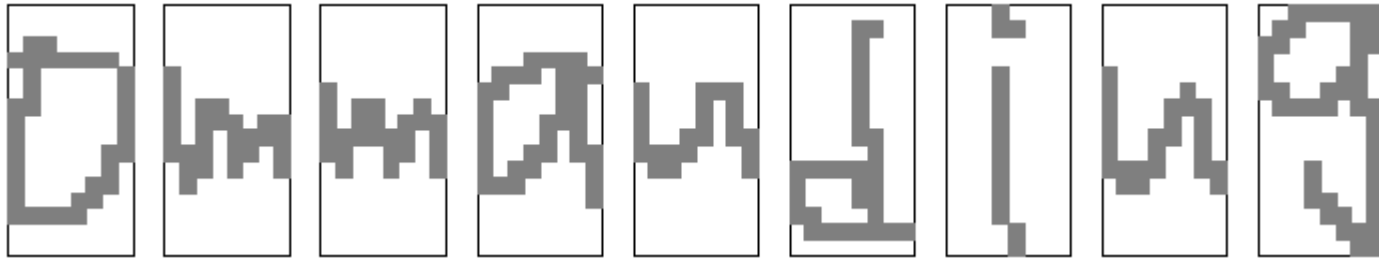## Graph

## Sequence of Classifications

# Structured Prediction

- Input to Classifier:
  - Local image features in neighborhood of pixel
  - Current neighboring pixels' labels

- Neighboring labels depend on classifier itself

- DAgger finds a classifier that does well at predicting pixel labels given the neighbors' labels it itself generates during the labeling process.

# Experiments: Handwriting Recognition
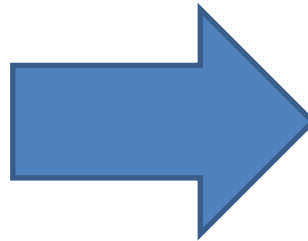


[Taskar 2003]

## Input:

Image current letter:



Previous predicted letter: O

## Output:

Current letter in {a,b,...,z}

# Test Folds Character Accuracy