# Validating Complex Agent Behavior

Scott Wallace

University of Michigan

University of Michigan – May 2003

# The Problem of Correctness



- Agents must have correct, expert-level behavior
- Errors undermine project's goals
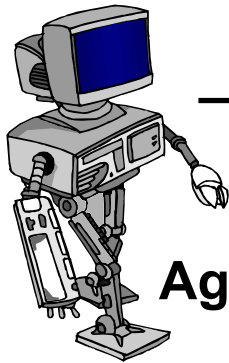- How can we ensure correctness?

# The Validation Bottleneck

- Manual Validation: Expert critiques agent behavior
  - Requires significant human effort
  - Difficult to detect every error
  - Standard approach to obtaining correct behavior

- Challenges for Automated Validation
  - *Difficult to formalize and articulate parameters of correct/incorrect behavior*
  - *"I can't tell you what's incorrect, but I know it when I see it."*
  - *Removing humans from the process creates new opportunities for failure*
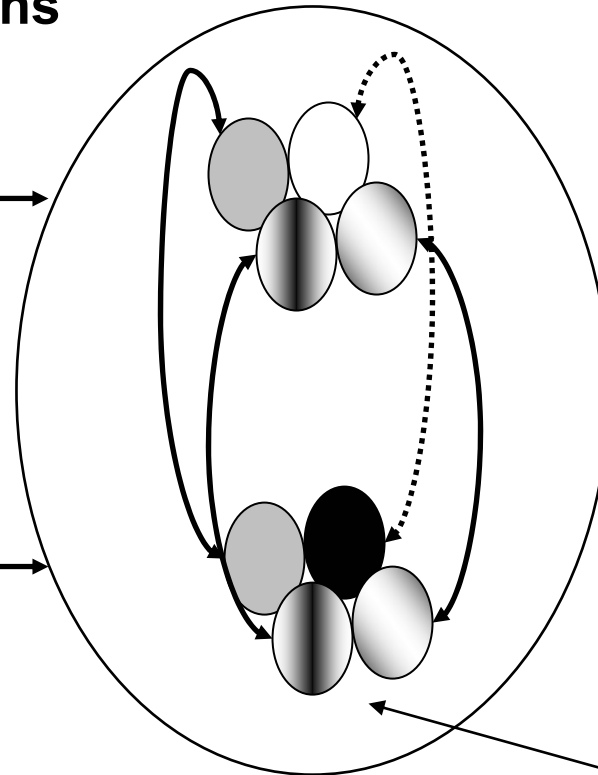
# Validation Framework Overview
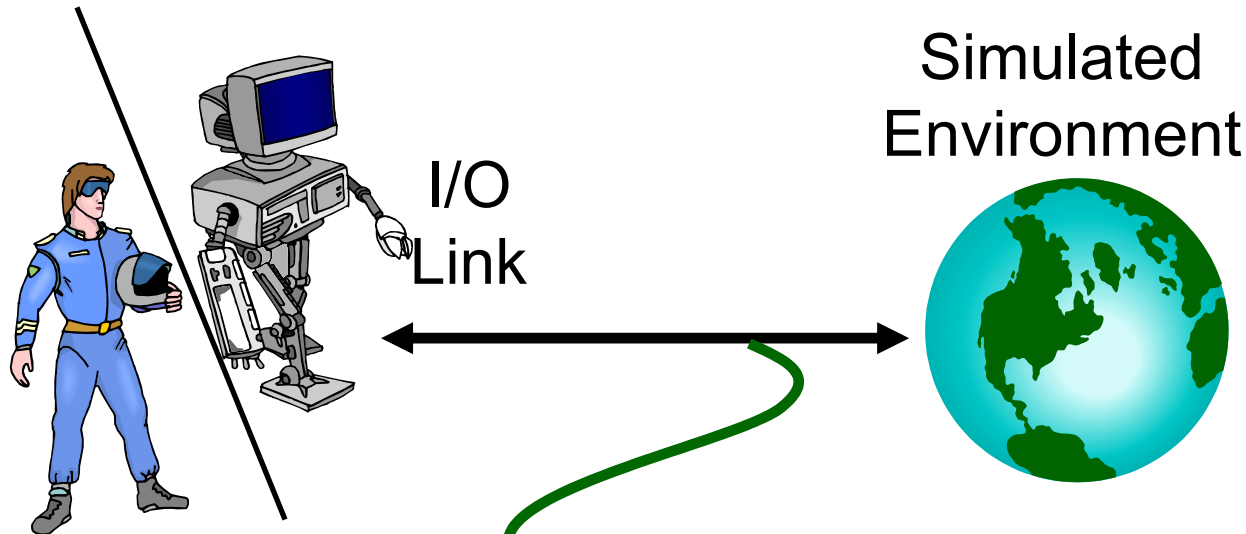


**Behavior Specifications**

Expert

Agent

**Comparison System**

**Summary of Behavior Differences**

Internal Behavior Representations

# Behavior Specifications



Simulated Environment
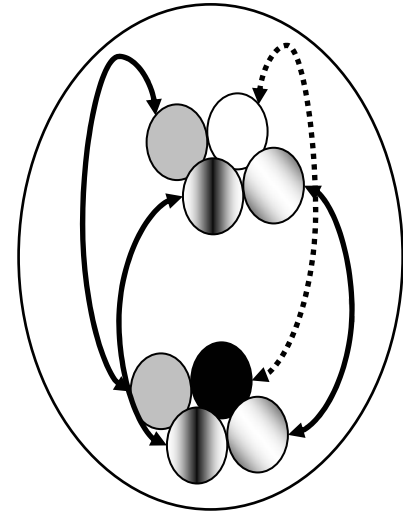
I/O Link

*Captured Actor I/O:* (s,a)

■ Actor interacts with simulator

■ Simulator provides a clean interface for:

   – Identifying salient state information

   – Identifying relevant actions

*Actor's Goal Annotations:* (s,G)

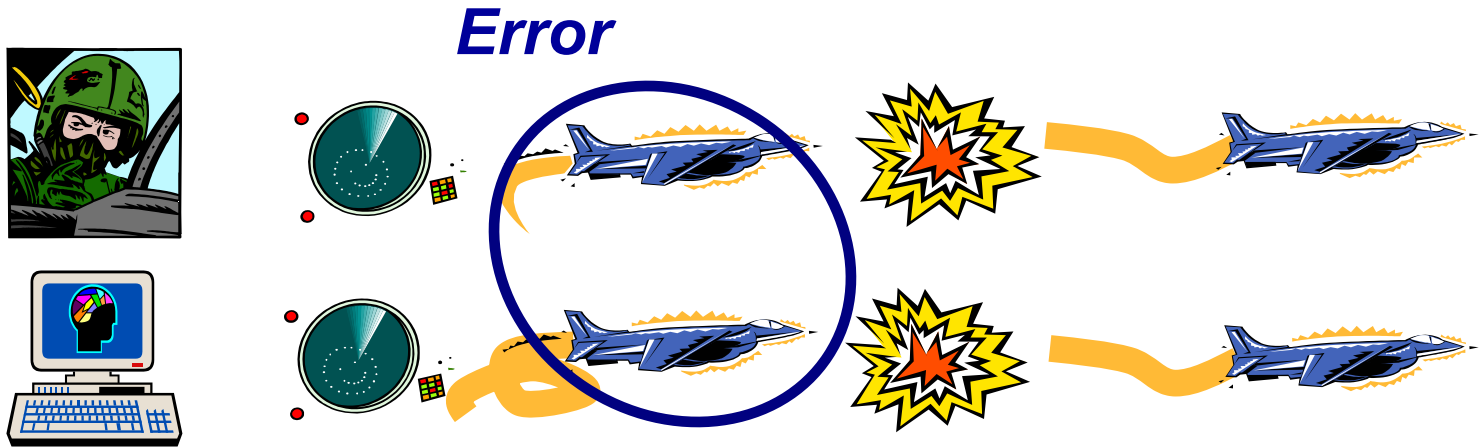$(s, G, a)_1, \ldots, (s, G, a)_n$

Behavior Trace

# Comparison System

- Desirable attributes:
  - Low Human/Computational Effort
  - Domain Independence
  - Efficacy

- We examine two types of approaches
  - Sequential (actions, goals)
  - Behavior bounding

- Quality of the comparison system will be influenced by choice of representation
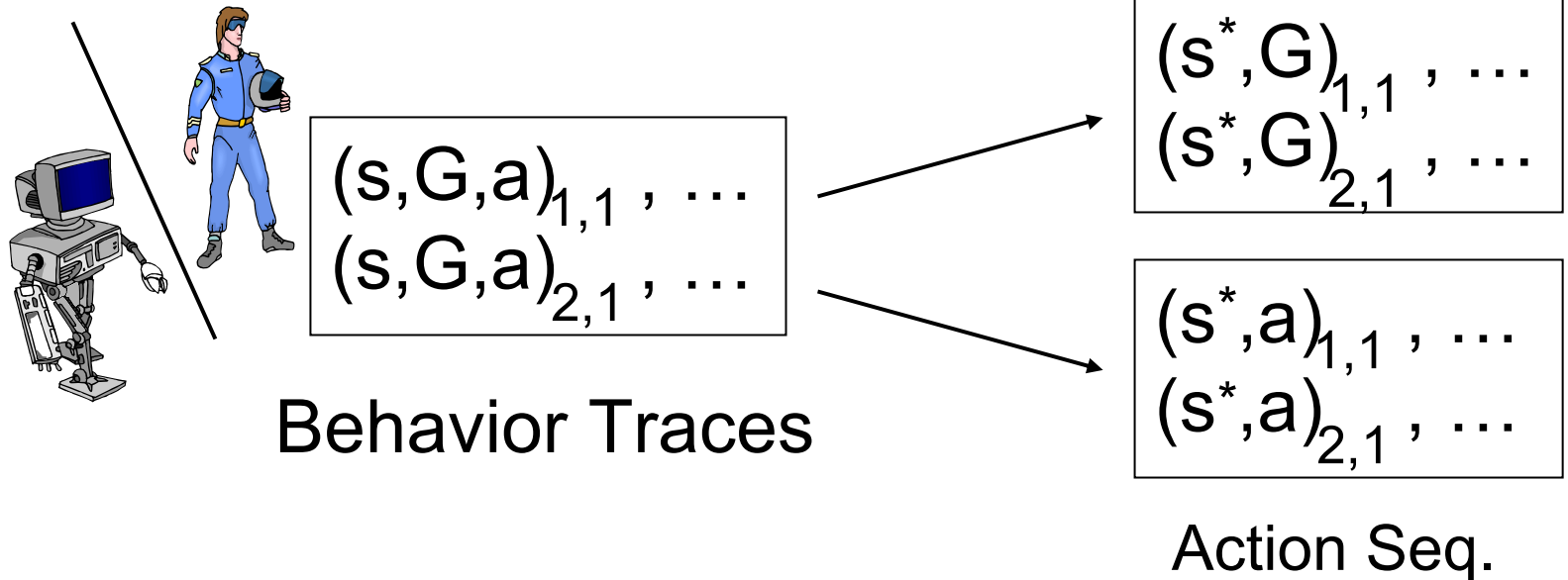
# Overview: Sequential Approach

*Error*



■ Discrepancies between sequences indicate errors

# Sequential Approaches



Goal Seq.

$(s^*,G)_{1,1}$ , …
$(s^*,G)_{2,1}$ , …

$(s,G,a)_{1,1}$ , …
$(s,G,a)_{2,1}$ , …

Behavior Traces

$(s^*,a)_{1,1}$ , …
$(s^*,a)_{2,1}$ , …

Action Seq.

- Extract symbols from behavior traces to form sequences (internal behavior representations)

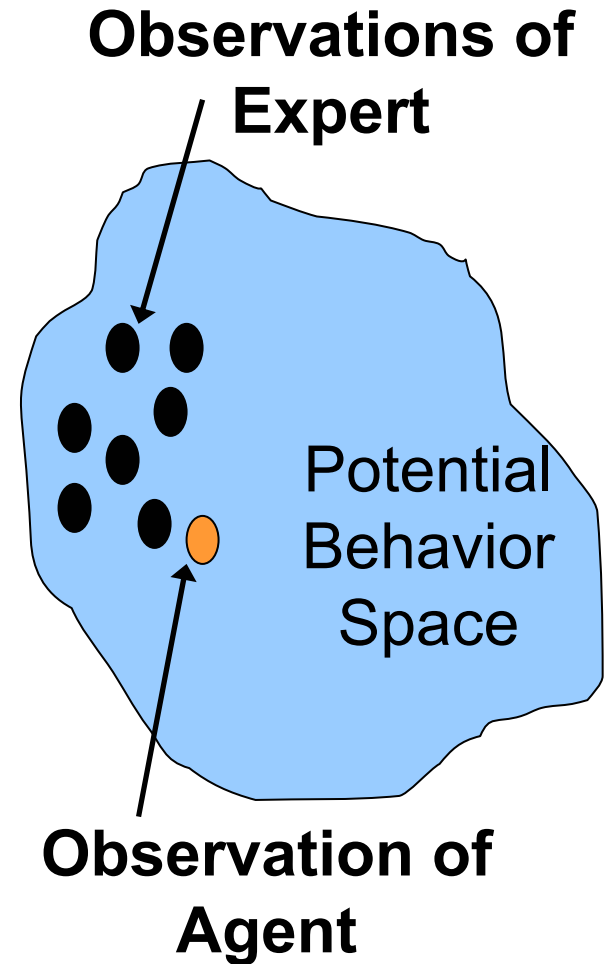- Compare sequences once aligned to minimize differences

# Effort and Domain Independence

- Sequences are very weak generalizations of behavior traces
- Required expert examples grows rapidly with:
  - Complexity of domain/behavior
  - Variability of behavior
- Internal representation grows with number of expert examples
- Computational complexity (time/space) of comparison is a function of representation size

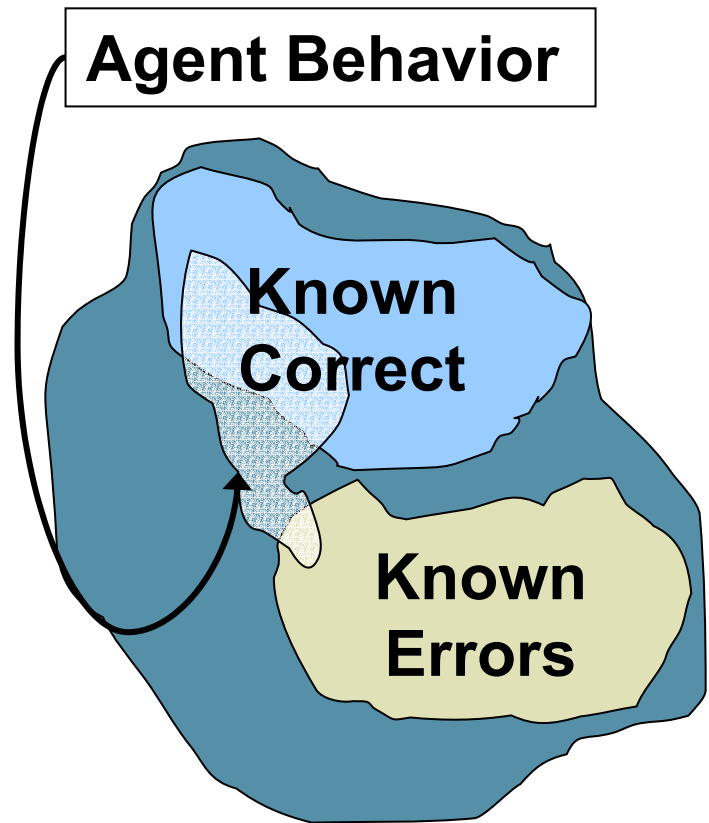- *But* representation makes few assumptions

# Weakness of Sequential Approach

- Sequences represent instances of behavior
- Instances are points in the behavior space
- Want to represent aggregate behavior

**Observations of Expert**

Potential Behavior Space
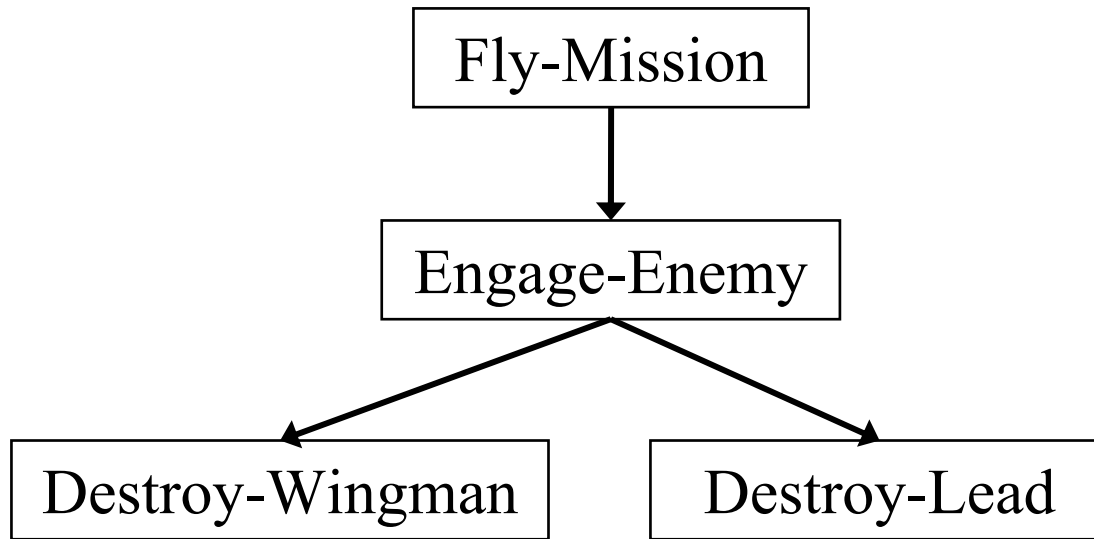
**Observation of Agent**

# Behavior Bounding

- Define boundaries in the space of potential behavior using:
  - observations
  - knowledge of task requirements
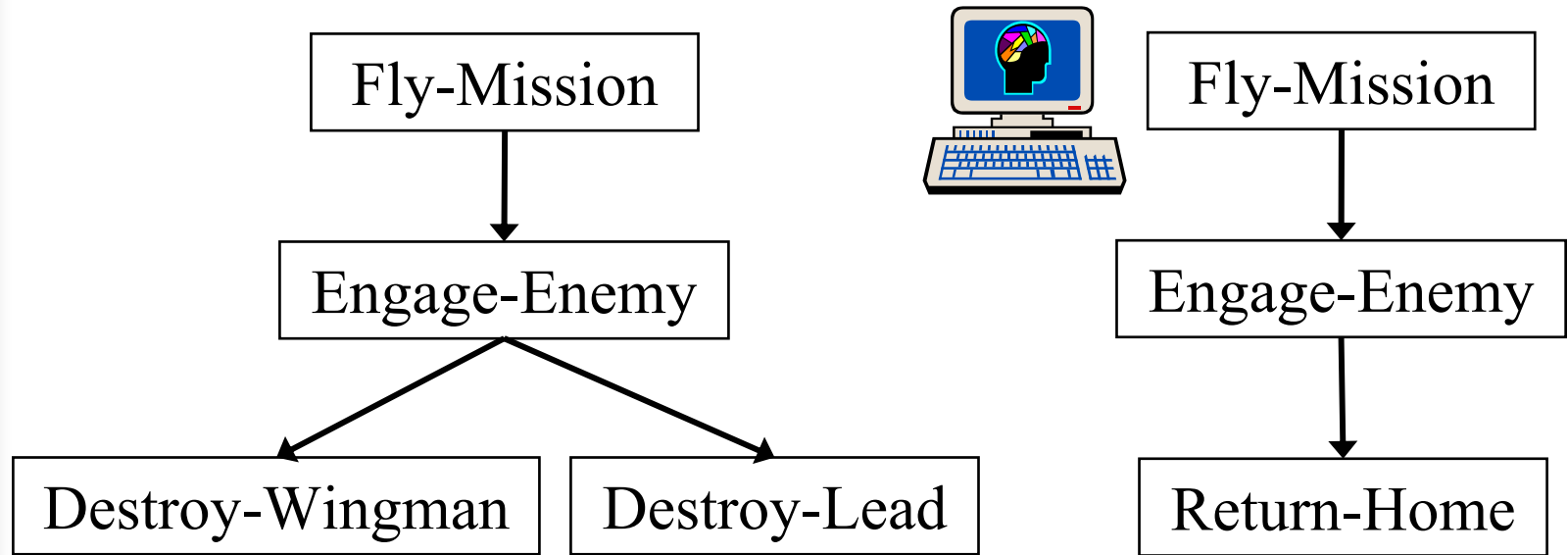- Determine portion of agent behavior in each region

**Agent Behavior**

**Known Correct**

**Known Errors**

# Leveraging the Goal Hierarchy

```
        ┌─────────────┐
        │ Fly-Mission │
        └─────────────┘
               │
               ▼
       ┌──────────────┐
       │ Engage-Enemy │
       └──────────────┘
          ╱        ╲
         ▼          ▼
┌──────────────────┐  ┌──────────────┐
│ Destroy-Wingman  │  │ Destroy-Lead │
└──────────────────┘  └──────────────┘
```
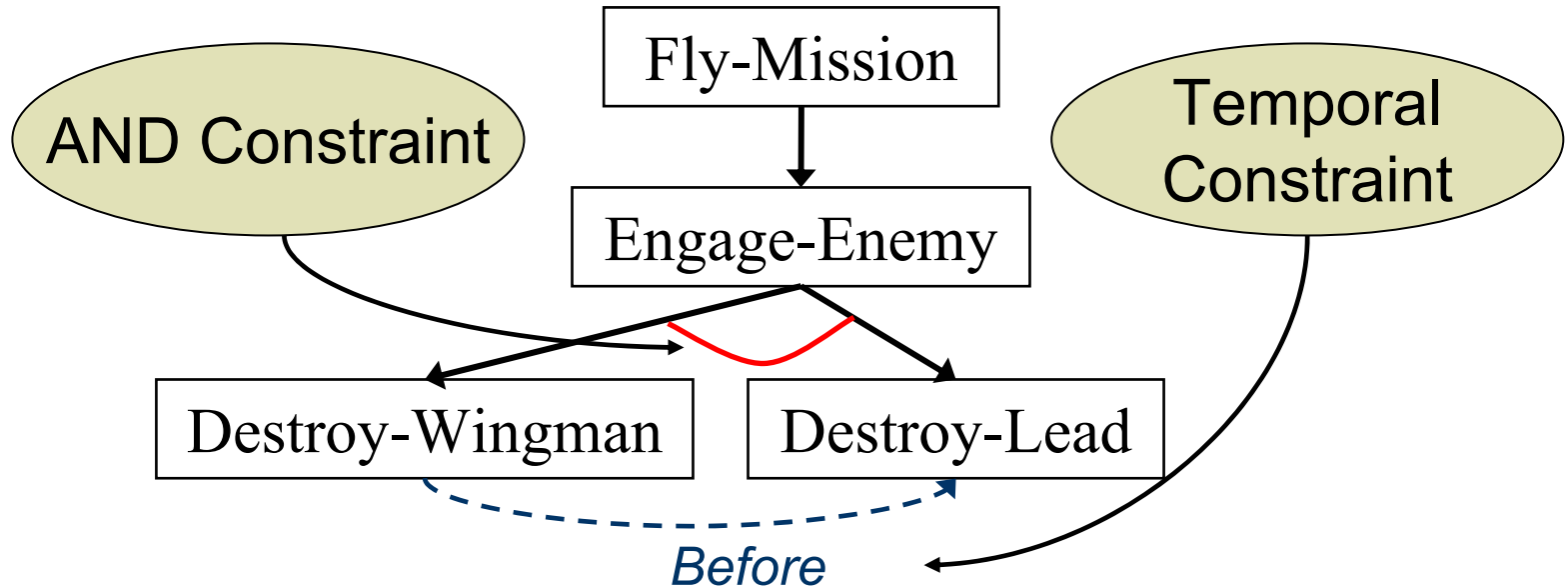
- A hierarchy compactly represents a subset of the behavior space.

- Agents are often constructed via task-decomposition.

- A hierarchy can be built from behavior traces.

# Goal Hierarchies as Classifiers

```
┌─────────────┐                      ┌─────────────┐
│ Fly-Mission │                      │ Fly-Mission │
└─────────────┘                      └─────────────┘
       │                                    │
       ▼                                    ▼
┌───────────────┐                    ┌───────────────┐
│ Engage-Enemy  │                    │ Engage-Enemy  │
└───────────────┘                    └───────────────┘
      ╱       ╲                             │
     ▼         ▼                            ▼
┌──────────────┐ ┌─────────────┐     ┌─────────────┐
│Destroy-Wingman│ │ Destroy-Lead│     │ Return-Home │
└──────────────┘ └─────────────┘     └─────────────┘
```

- Basic hierarchy identifies differences in topology.
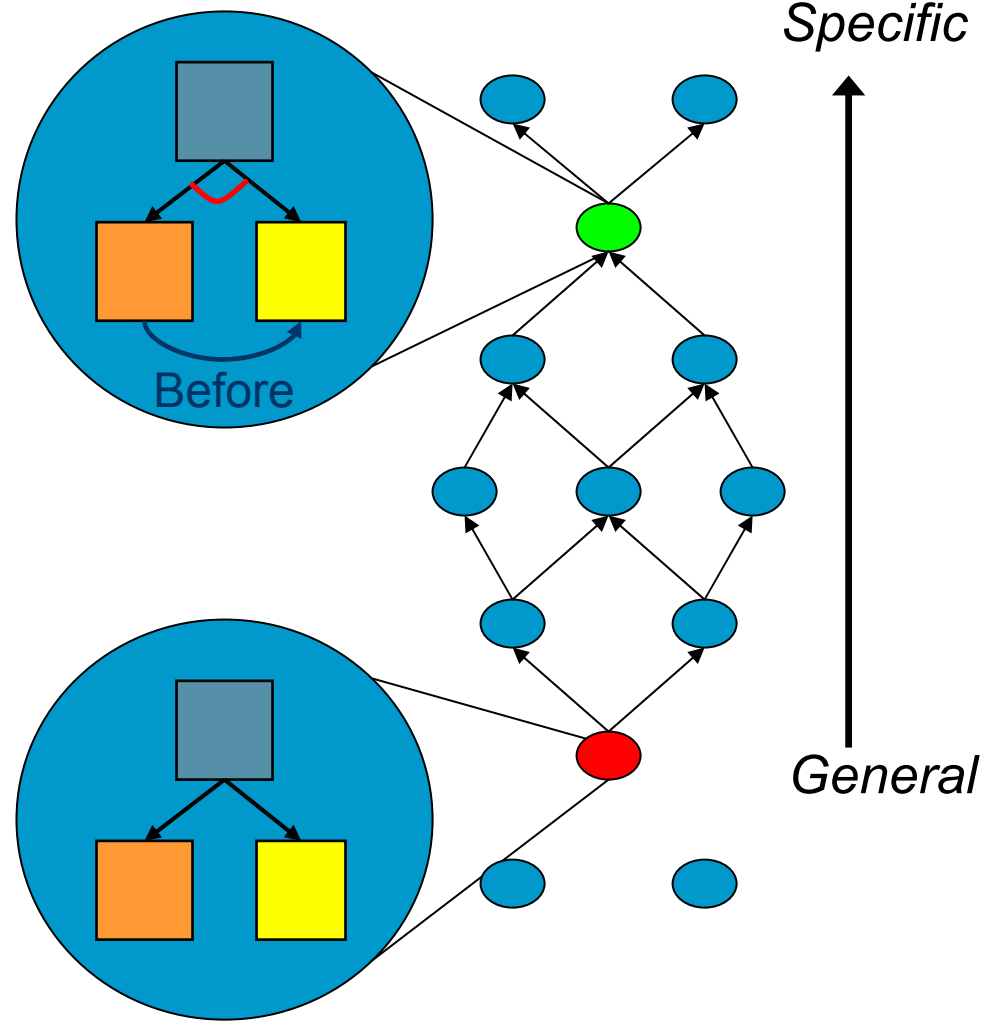
# Constrained Goal Hierarchy



- Constraints reduce degrees of freedom
- Create specializations of original hierarchy
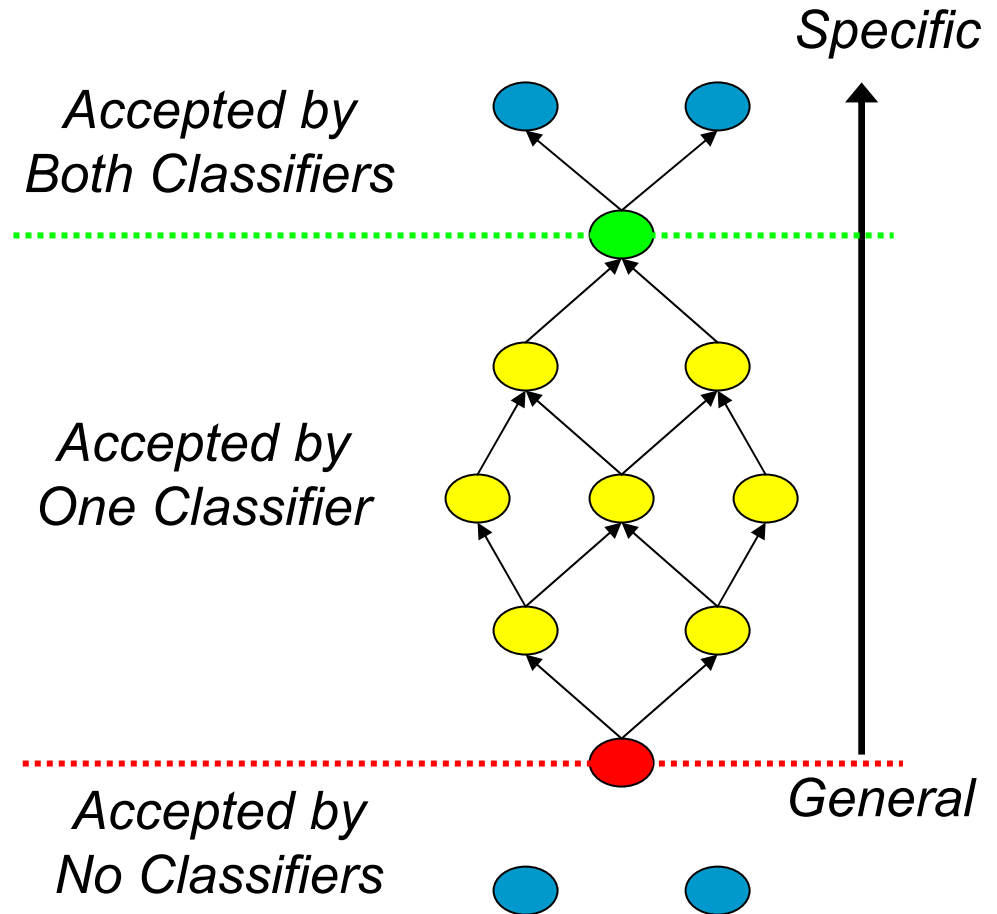- Can also be used to classify behavior

# Hierarchies As Partitions

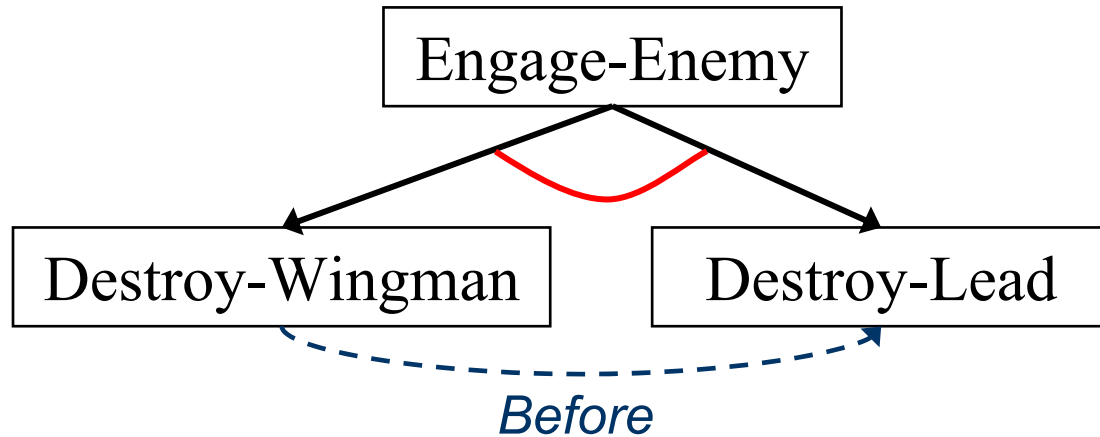- Constraints impose an ordering on the behavior space

# Hierarchies As Partitions

- Boundaries partition space into three regions

- Paves way for efficient error detection

*Specific*

*Accepted by Both Classifiers*

*Accepted by One Classifier*

*Accepted by No Classifiers*

*General*

# Building an Upper Boundary

Engage-Enemy

Destroy-Wingman          Destroy-Lead

*Before*

*Fly-Mission*          *Fly-Mission*
*Engage-Enemy*          *Engage-Enemy*
*Destroy-Wingman*          *Destroy-Lead*

- Construct a maximally specific hierarchy covering the observations

# Building a Lower Boundary

```
        ┌─────────────┐
        │ Fly-Mission │
        └─────────────┘
               │
               ▼
       ┌───────────────┐
       │ Engage-Enemy  │
       └───────────────┘
          ╱          ╲
         ▼            ▼
┌──────────────────┐  ┌──────────────┐
│ Destroy-Wingman  │  │ Destroy-Lead │
└──────────────────┘  └──────────────┘
```
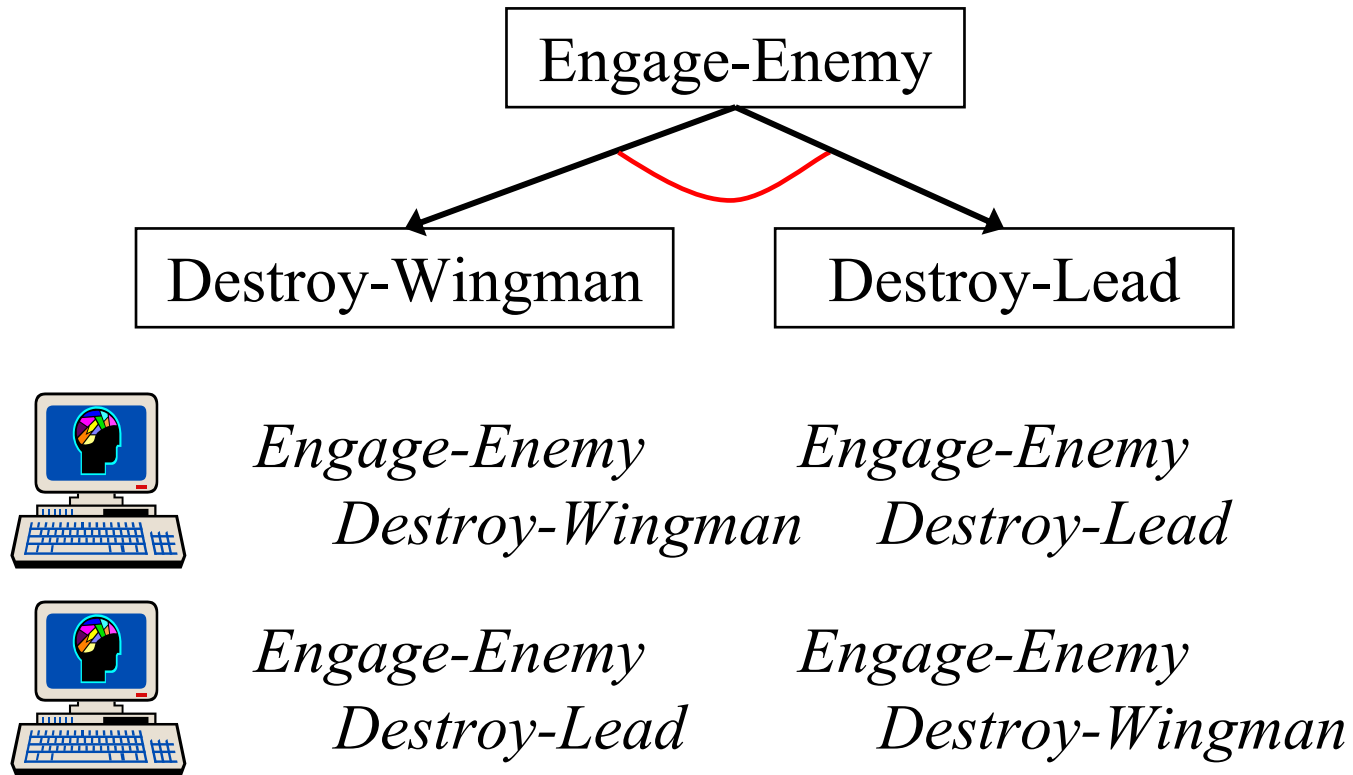
- Removing constraints yields lower bound
- Alternatively, lower bound may be generated manually.

# Partition Behavior Space

# Observe Agent Behavior
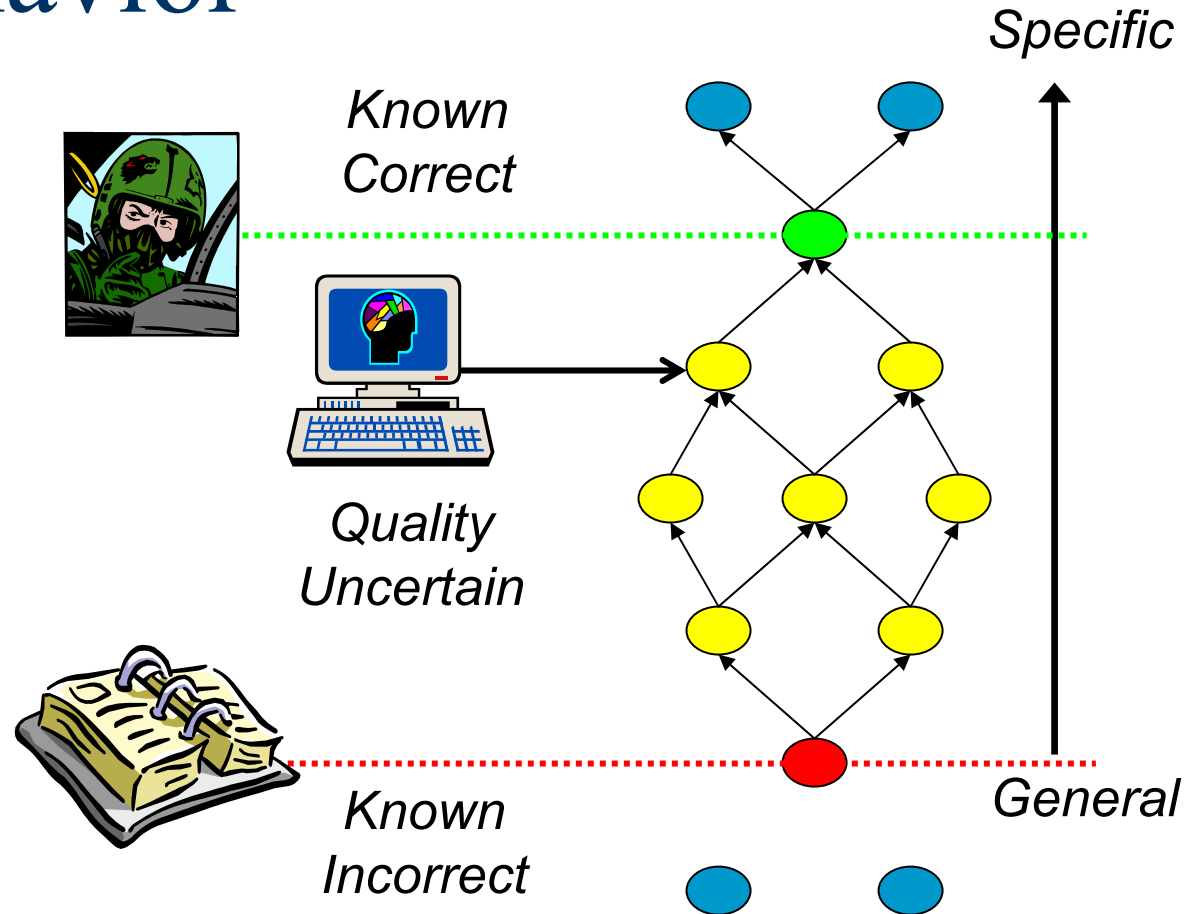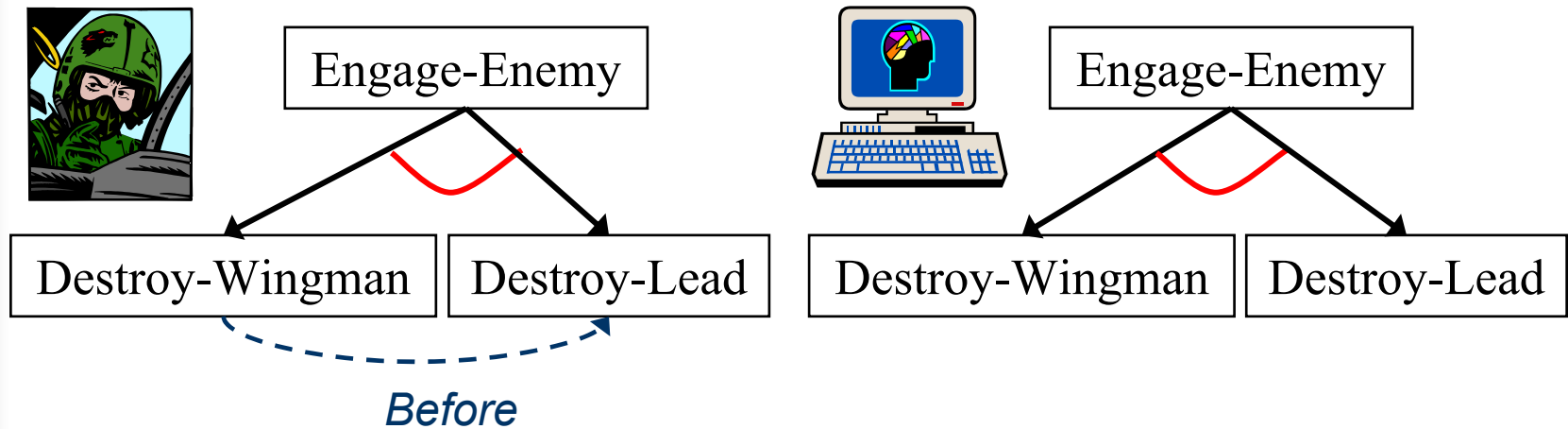
```
                    ┌──────────────────┐
                    │  Engage-Enemy    │
                    └──────────────────┘
                      ╱              ╲
                     ╱                ╲
        ┌────────────────────┐  ┌──────────────────┐
        │  Destroy-Wingman   │  │  Destroy-Lead    │
        └────────────────────┘  └──────────────────┘
```

*Engage-Enemy*          *Engage-Enemy*
  *Destroy-Wingman*        *Destroy-Lead*

*Engage-Enemy*          *Engage-Enemy*
  *Destroy-Lead*           *Destroy-Wingman*

- Construct a maximally specific hierarchy covering the observations

# Identify Quality of Agent Behavior

# Identify Quality of Agent Behavior



Engage-Enemy

Destroy-Wingman    Destroy-Lead

*Before*

Engage-Enemy

Destroy-Wingman    Destroy-Lead

- Agent behavior is not a specialization of Expert behavior
- Looking at behaviors encapsulated by hierarchy gives details of similarities and differences
  - Agent may perform sub-goals in an incorrect order

# Effort and Domain Independence

- Hierarchies can be built using relatively few behavior traces

- Computation effort of comparison
  - Independent of number of expert examples
  - Polynomial in size of hierarchy

- *Representation should be compatible with many goal based agents*

# Measuring Efficacy

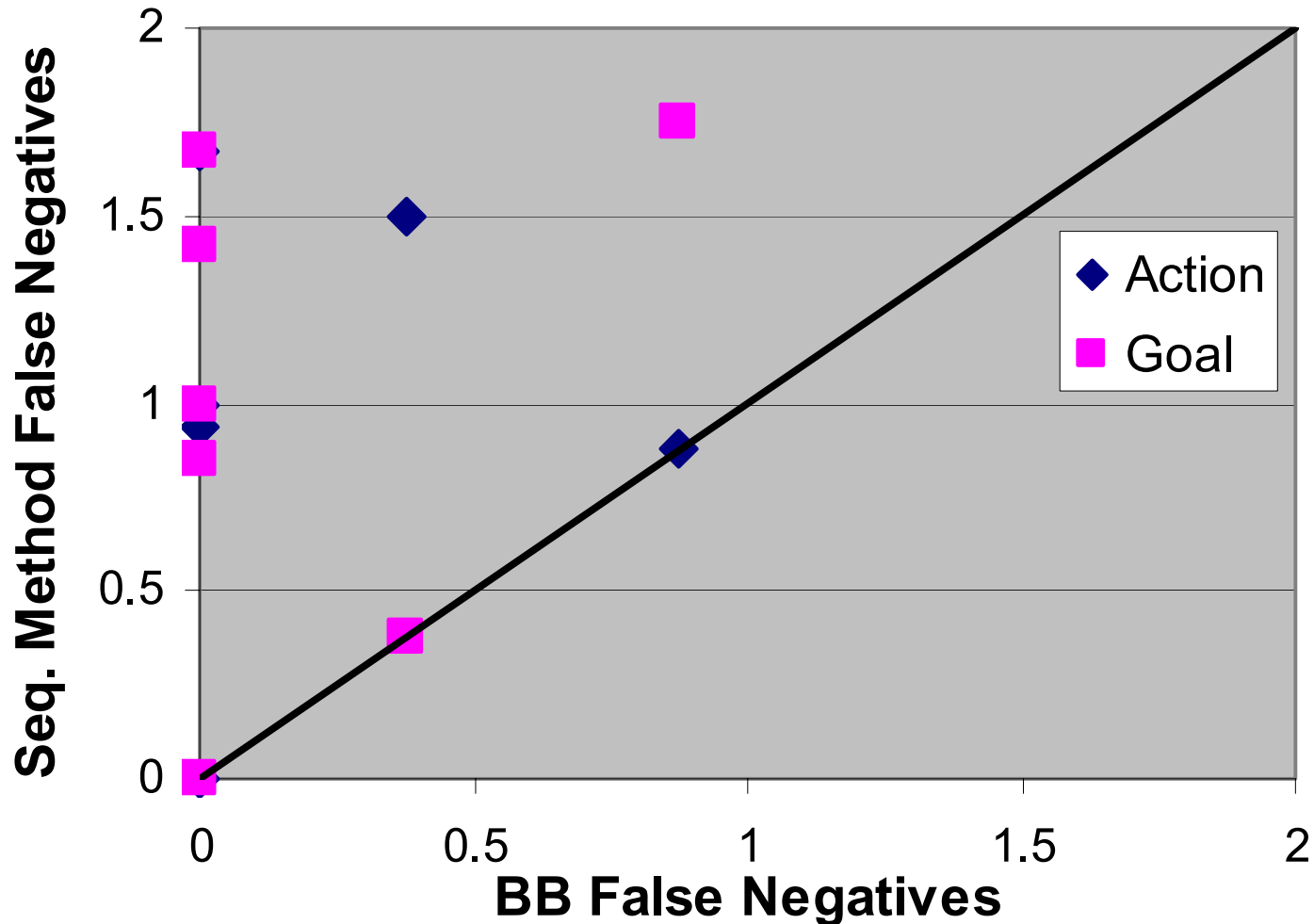- Rate method based on the quality of data in the summary

$$Report\ Density = \frac{True\ Positives}{Reported}$$
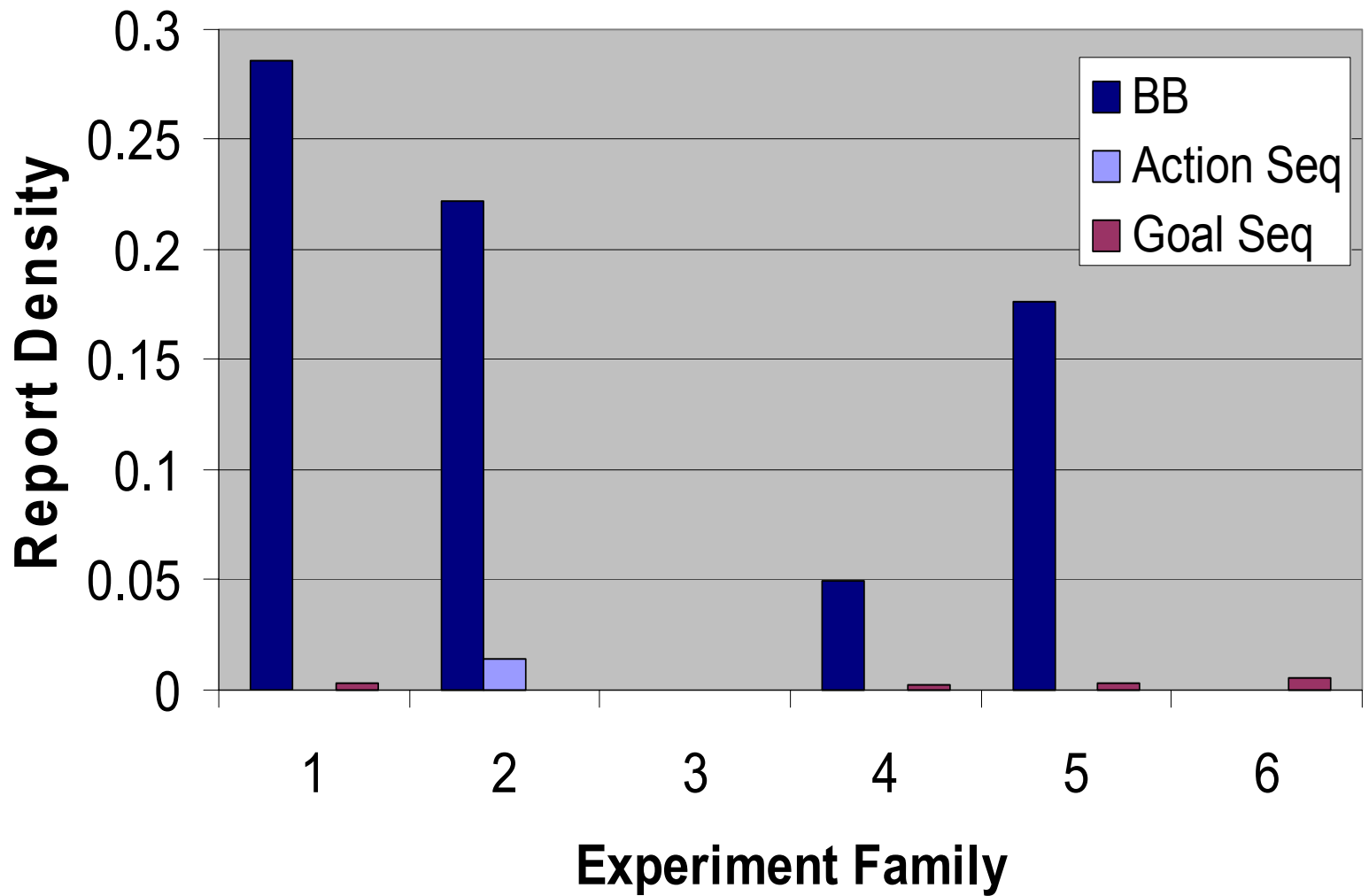
- Also want few undetected errors: *False Negatives*
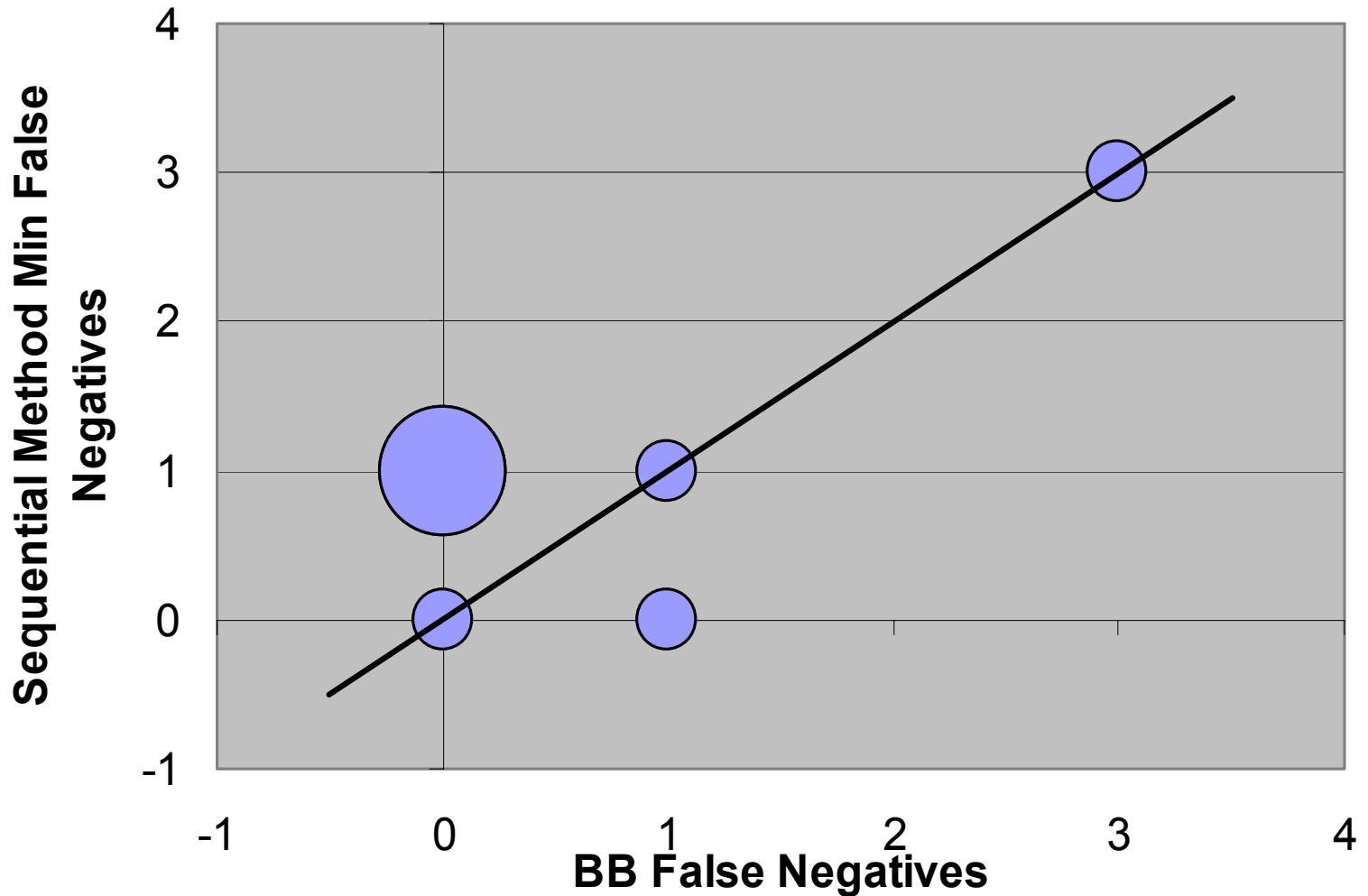
# Results: Object Retrieval Domain

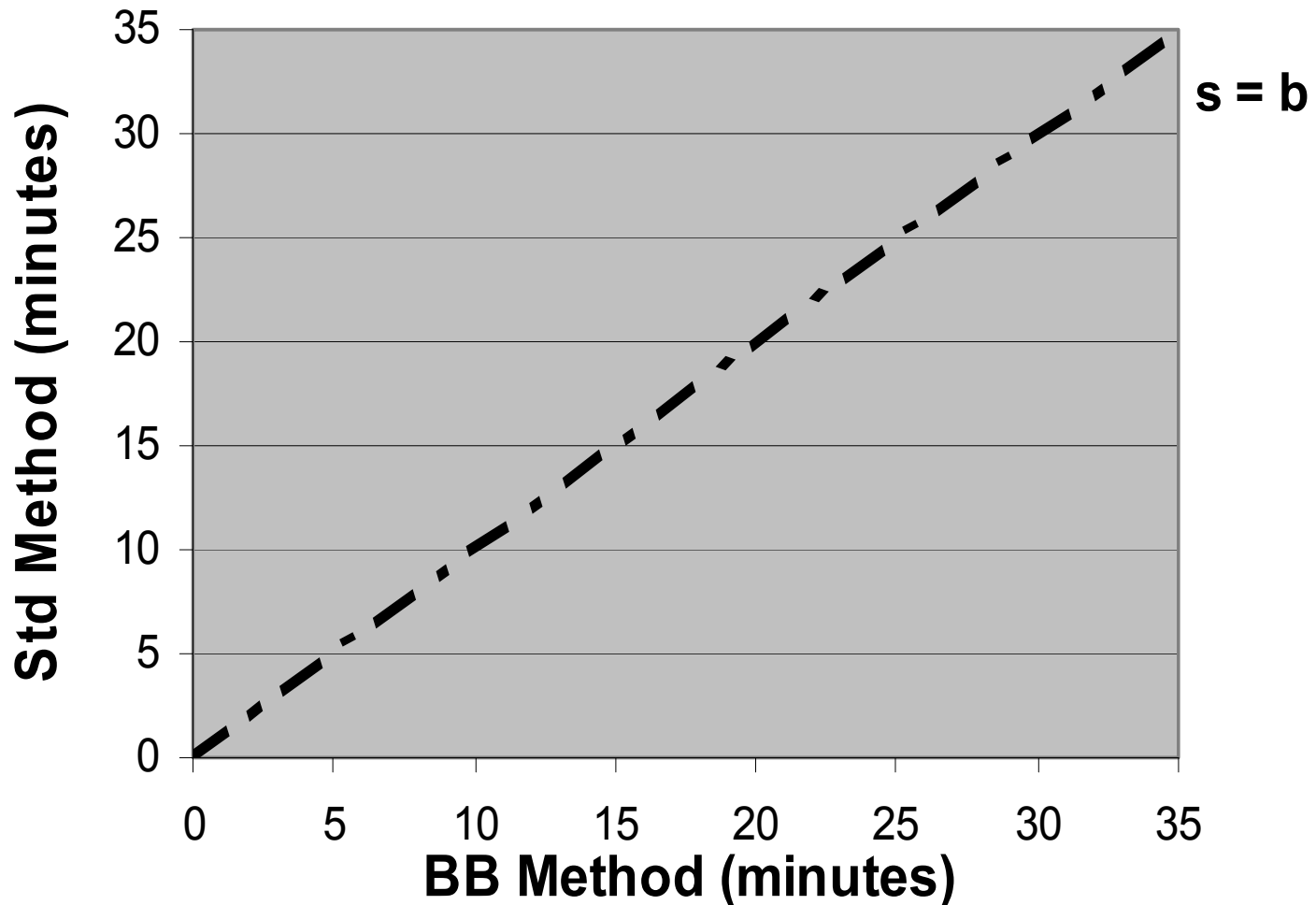# False Negatives in Object Retrieval Domain
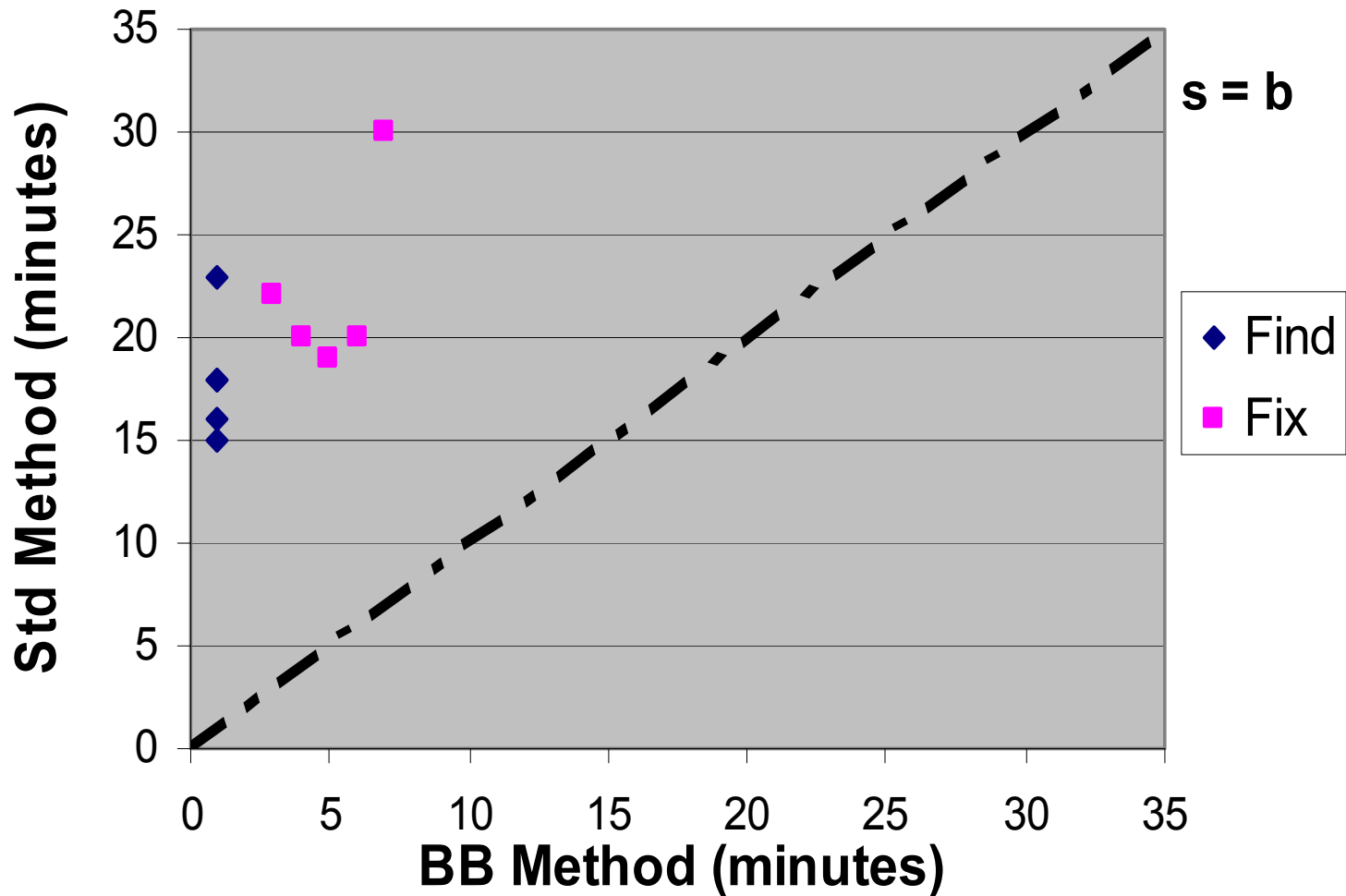
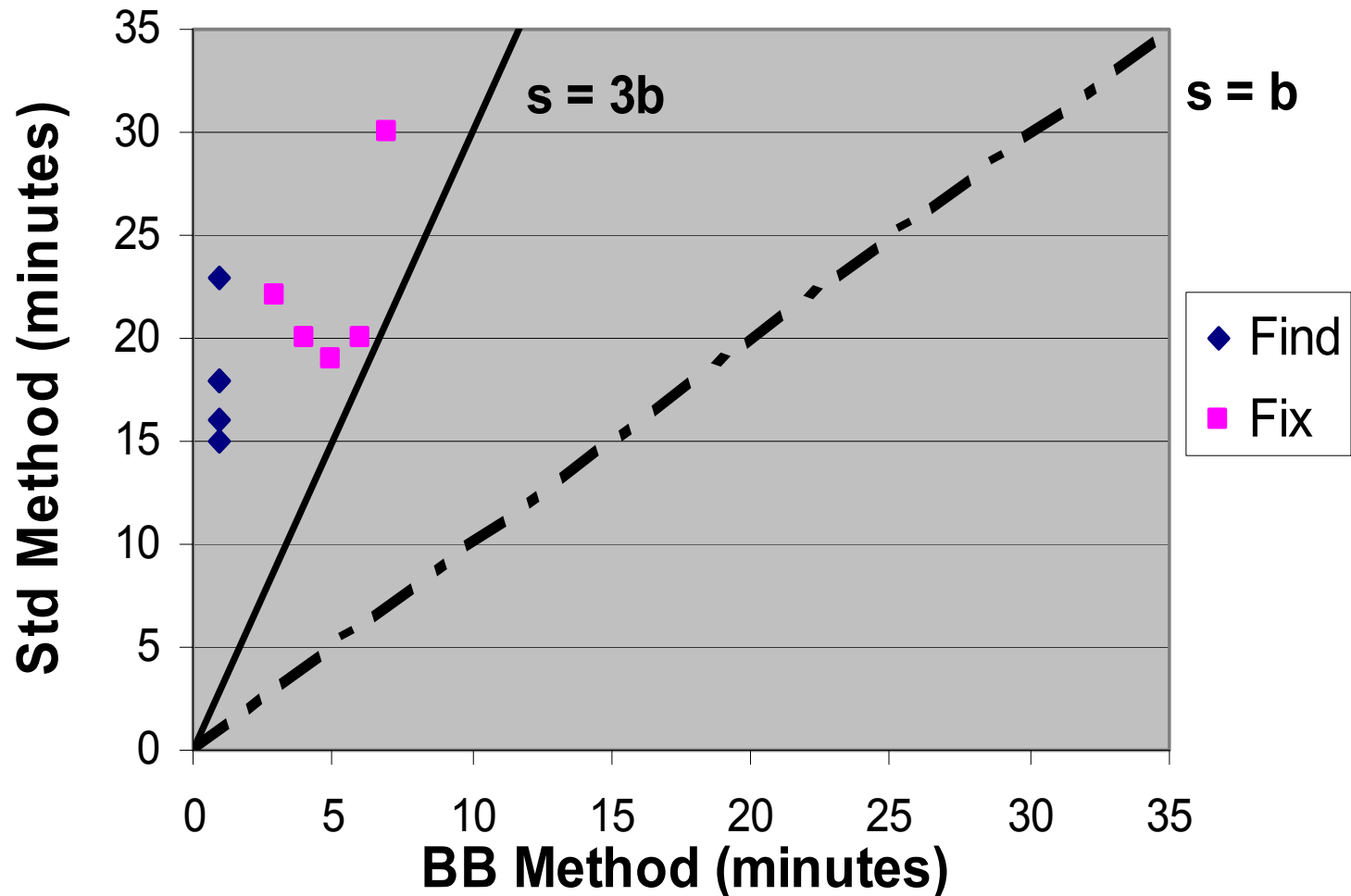# Results: MOUT Domain

# False Negatives in MOUT Domain

# Behavior Bounding as a Validation Tool

# Behavior Bounding as a Validation Tool

# Behavior Bounding as a Validation Tool

# Nuggets & Coal

- Simple, general behavior representation
- Leverages the natural organization of Soar agents
- Low cost to generate
- Performs well compared to sequential approach

- Simplicity leaves it susceptible to overgeneralization