# SOS: Safe, Optimal* and Small Strategies for Hybrid Markov Decision Processes

Pranav Ashok[1], Jan Kretinsky[1], Kim Guldstrand Larsen[2], Adrien Le Coënt[2], Jakob Haahr Taankvist[2] and Maximilian Weininger[1]

[1]Technical University of Munich, Munich, Germany
[2]Aalborg University, Aalborg, Denmark

# Outline

1. UPPAAL Stratego

2. Safe strategies

3. Compression of strategies

4. Our proposal: Stratego+

5. Results

6. Future directions

# What is Stratego?

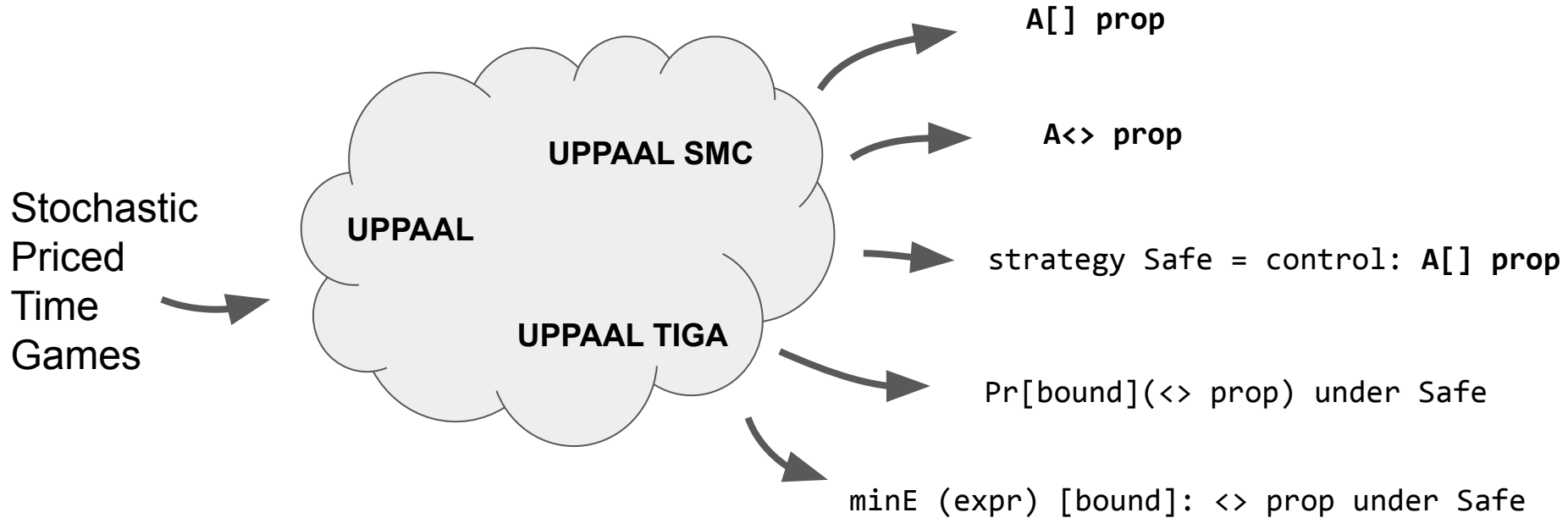Generate, optimize, evaluate, compare strategies

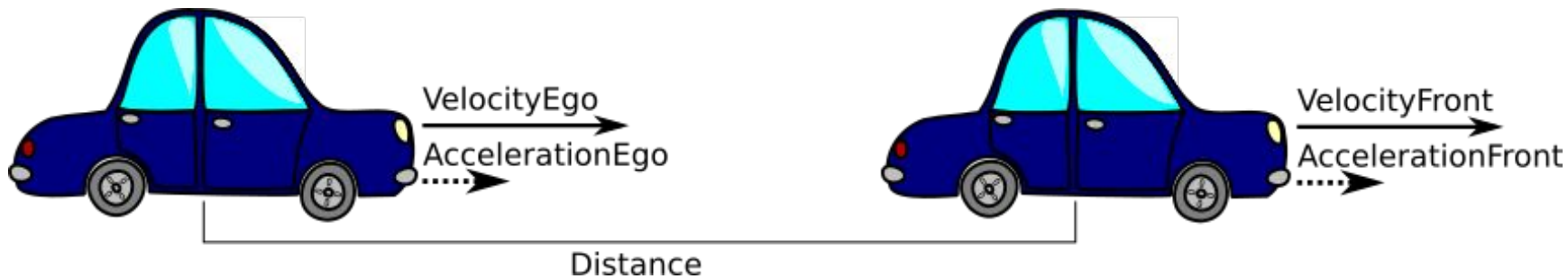**Hybrid Markov Decision Process**
costs, time, data variables ...

# What is Stratego?

The tool allows for ==efficient and flexible "strategy-space" exploration before adaptation in a final implementation== by maintaining strategies as first class objects in the model-checking query language.
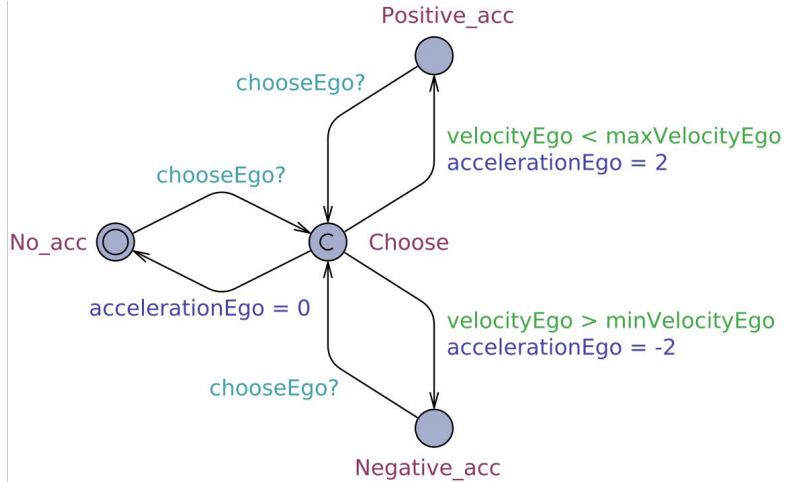
# What is Stratego?



Stochastic Priced Time Games → UPPAAL / UPPAAL SMC / UPPAAL TIGA →

- `A[] prop`
- `A<> prop`
- `strategy Safe = control: A[] prop`
- `Pr[bound](<> prop) under Safe`
- `minE (expr) [bound]: <> prop under Safe`

Strategies from UPPAAL Stratego are **huge** and **incomprehensible**

safety: A[] distance ≥ 5
optimality: minimize aggregate distance

Ego

Positive_acc

chooseEgo?

velocityEgo < maxVelocityEgo
accelerationEgo = 2

chooseEgo?

No_acc        C     Choose

accelerationEgo = 0

Front

velocityEgo > minVelocityEgo
accelerationEgo = -2

chooseEgo?

Negative_acc

safety: A[] distance ≥ 5
optimality: minimize aggregate distance

# Sample safe strategy

```
State: ( Ego.No_acc Front.No_acceleration ) distance=200 velocityEgo=16 accelerationEgo=0 velocityFront=12
accelerationFront=0
Wait.

State: ( Ego.Choose Front.Negative_acc ) distance=101 velocityEgo=-6 accelerationEgo=0 velocityFront=-4
accelerationFront=-2
Take transition Ego.Choose->Ego.No_acc
Take transition Ego.Choose->Ego.Positive_acc
Take transition Ego.Choose->Ego.Negative_acc

State: ( Ego.No_acc Front.Positive_acc ) distance=82 velocityEgo=2 accelerationEgo=0 velocityFront=10
accelerationFront=2
Wait.

State: ( Ego.No_acc Front.Positive_acc System.FrontNext Monitor._id12 ) distance=85
```
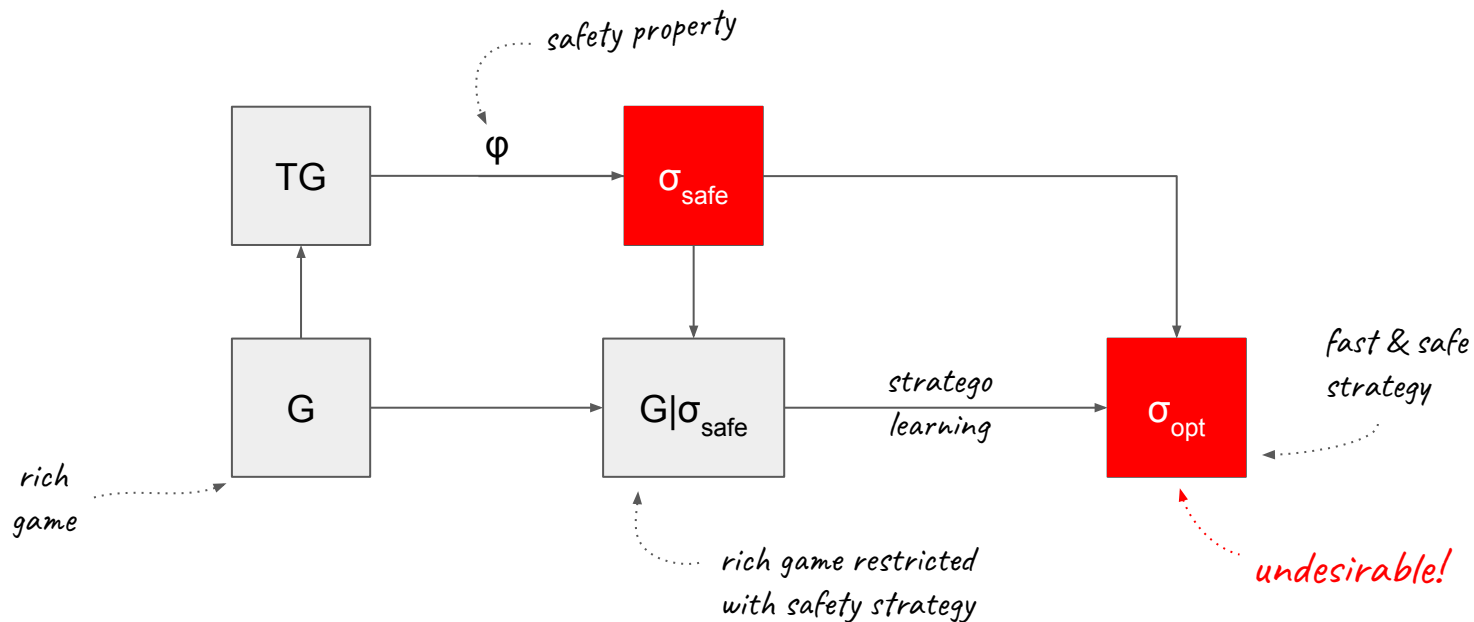
# Stratego Internals

Generating a safe and optimal controller for cruise control

# Problems

- 900k state action pairs, 300k controllable

- Incomprehensible for humans

- Lookup table too big for microcontrollers

- Executing inside Stratego takes too long

Would there exist a **small and safe** strategy?

# Can we learn a smaller representation?

# Can we learn a smaller representation?

Binary Decision Diagrams?

# Can we learn a smaller representation?

Each state is composed of variables with integer domains

```
State: ( Ego.No_acc Front.No_acceleration ) distance=200 velocityEgo=16 accelerationEgo=0 velocityFront=12
accelerationFront=0
Wait.

State: ( Ego.Choose Front.Negative_acc ) distance=101 velocityEgo=-6 accelerationEgo=0 velocityFront=-4
accelerationFront=-2
Take transition Ego.Choose->Ego.No_acc
Take transition Ego.Choose->Ego.Positive_acc
Take transition Ego.Choose->Ego.Negative_acc

State: ( Ego.No_acc Front.Positive_acc ) distance=82 velocityEgo=2 accelerationEgo=0 velocityFront=10
accelerationFront=2
Wait.

State: ( Ego.No_acc Front.Positive_acc System.FrontNext Monitor._id12 ) distance=85
```
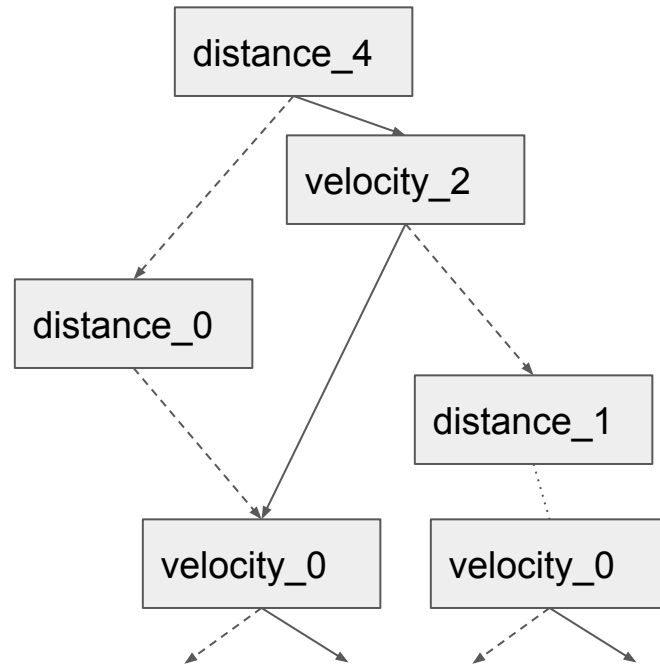
# Can we learn a smaller representation?

Each state is composed of variables with integer domains

# Can we learn a smaller representation?

Nearby points behave same?

```
State: ( Ego.No_acc Front.No_acceleration ) distance=200 velocityEgo=16 accelerationEgo=0 velocityFront=12
accelerationFront=0
Wait.

State: ( Ego.Choose Front.Negative_acc ) distance=101 velocityEgo=-6 accelerationEgo=0 velocityFront=-4
accelerationFront=-2
Take transition Ego.Choose->Ego.No_acc
Take transition Ego.Choose->Ego.Positive_acc
Take transition Ego.Choose->Ego.Negative_acc

State: ( Ego.No_acc Front.Positive_acc ) distance=82 velocityEgo=2 accelerationEgo=0 velocityFront=10
accelerationFront=2
Wait.

State: ( Ego.No_acc Front.Positive_acc System.FrontNext Monitor._id12 ) distance=85
```
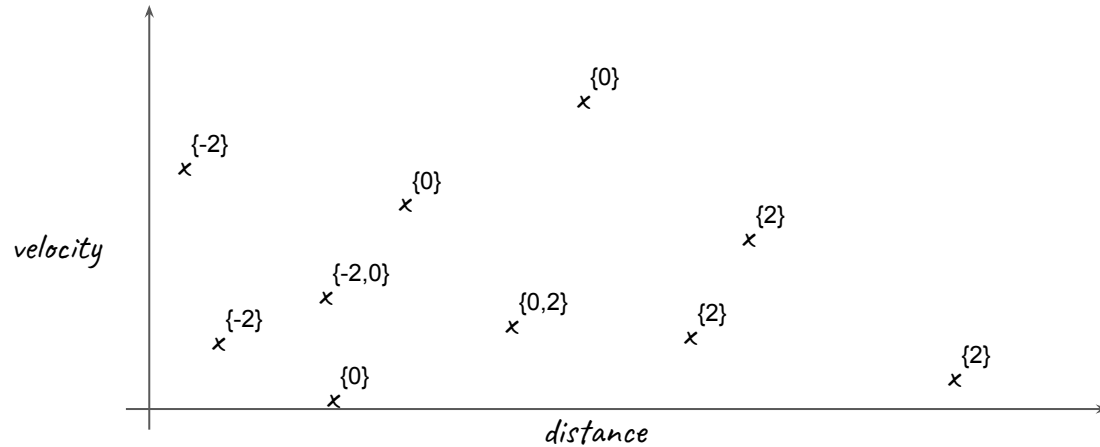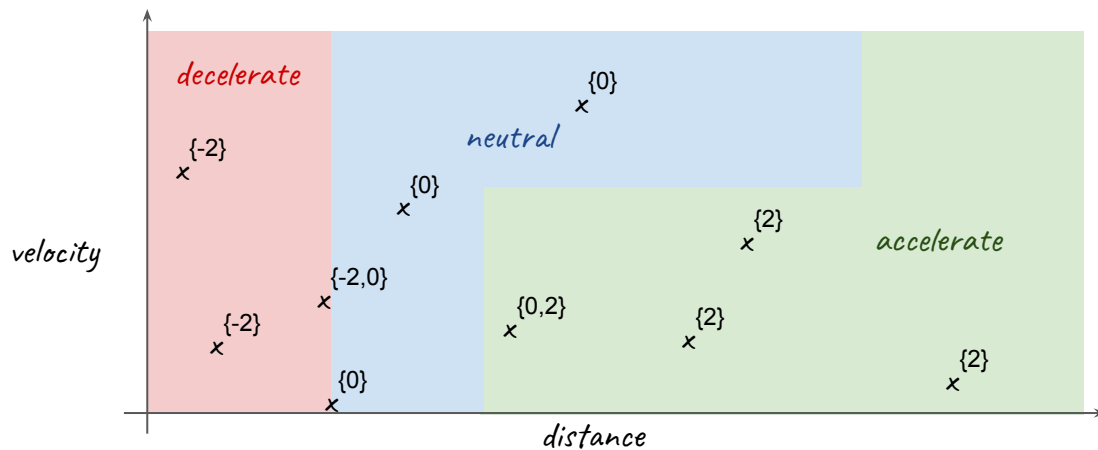
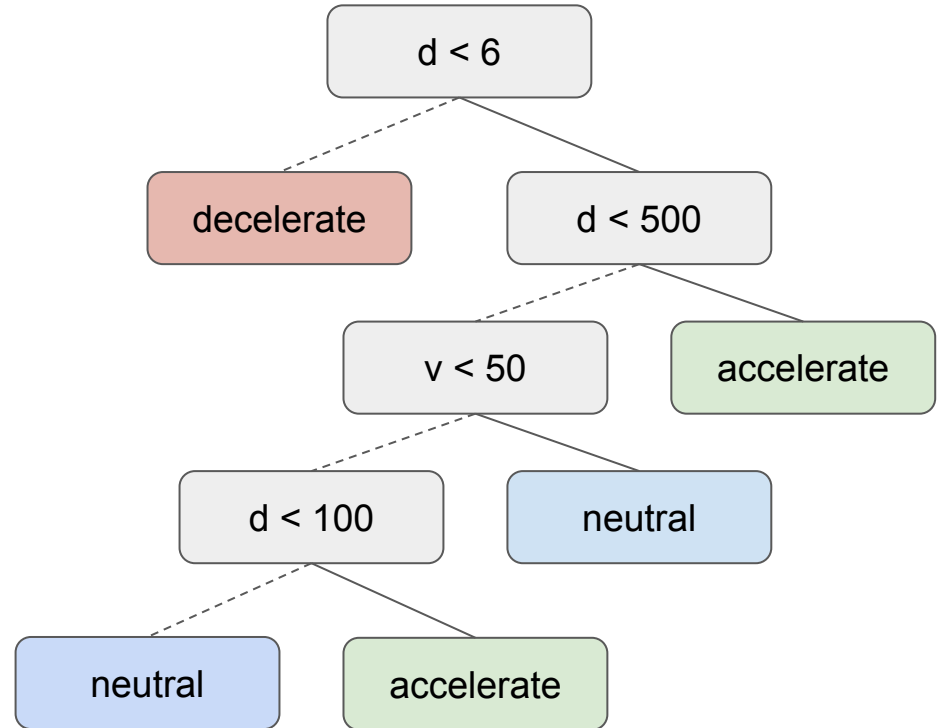# Can we learn a smaller representation?
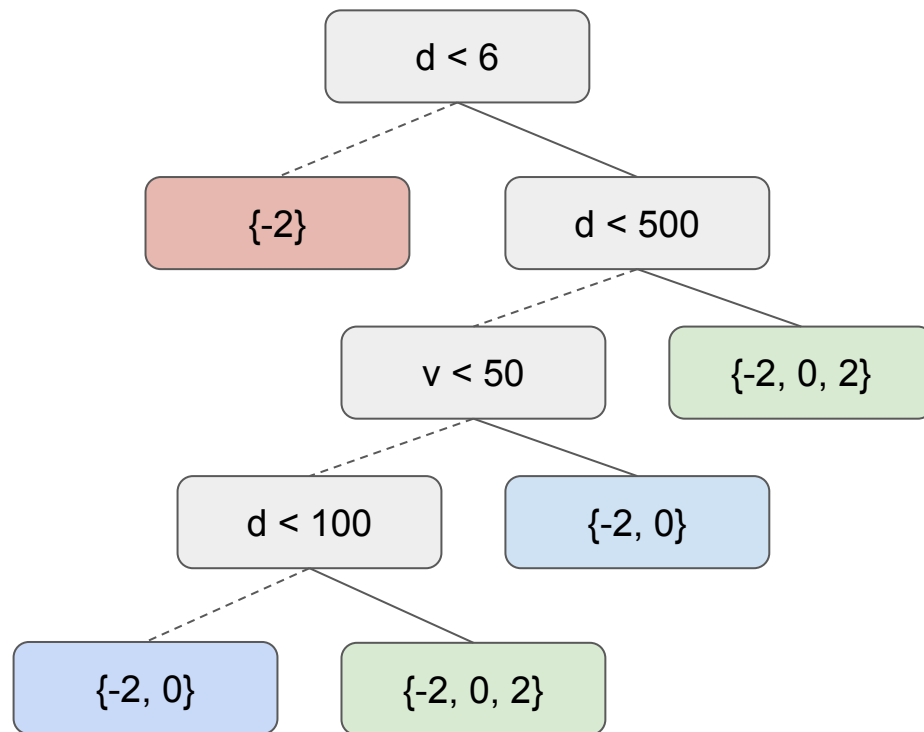
# Can we learn a smaller representation?

# Decision Trees

- No magic inside, unlike many other ML techniques

- Simple to interpret

- Good with features which have inherent ordering
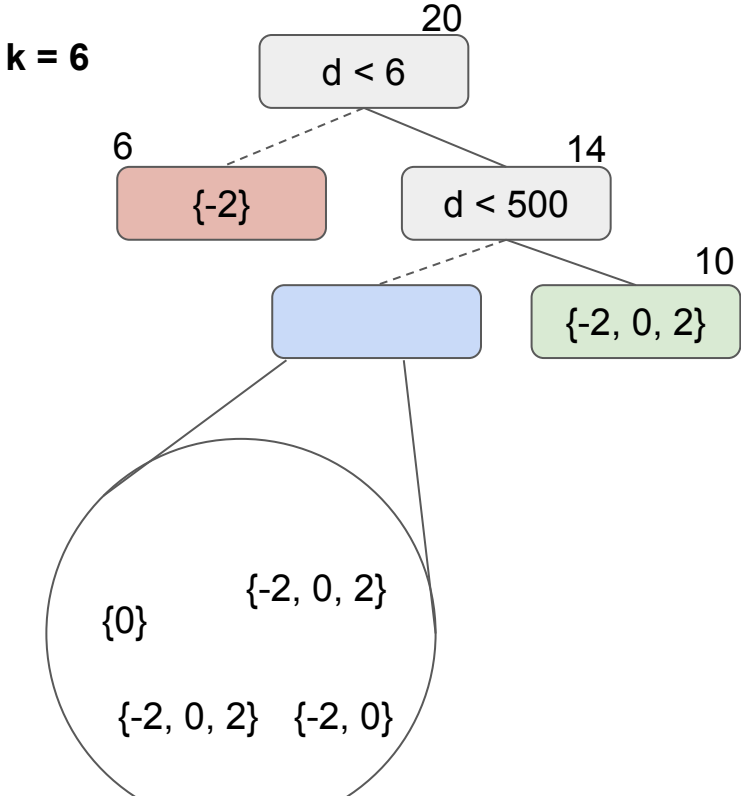
- Can be converted into executable code easily

# Construction of DTs

1. Multi-label classification

2. All leaves homogeneous
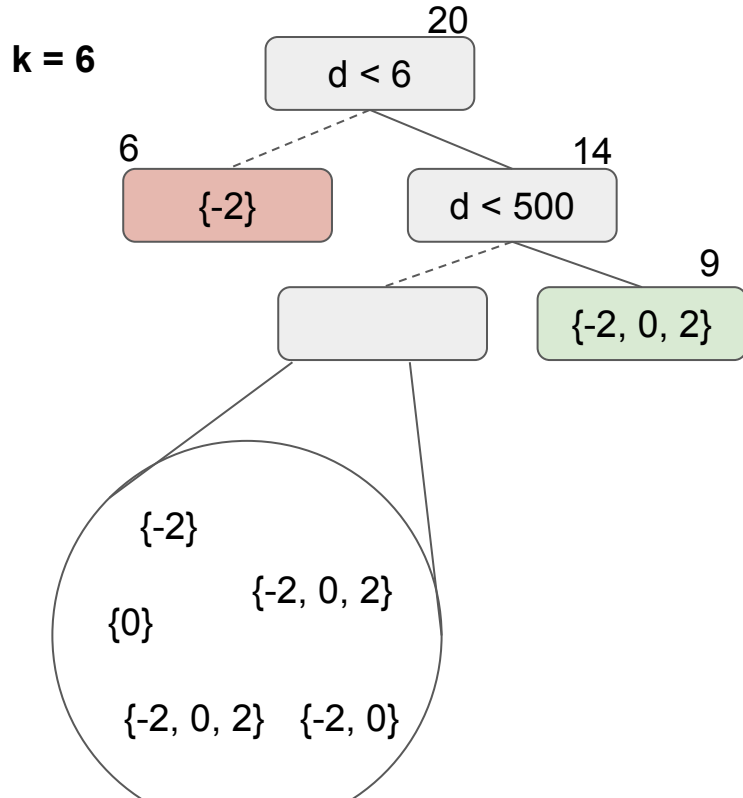
3. Size-permissiveness tradeoff

# Minimum split size

**k = 6**

20
d < 6

6
{-2}

14
d < 500

{-2, 0, 2}
10

{0}
{-2, 0, 2}
{-2, 0, 2}   {-2, 0}

Consider splitting node only if # data points > k

# Minimum split size

**k = 6**

20

d < 6

6

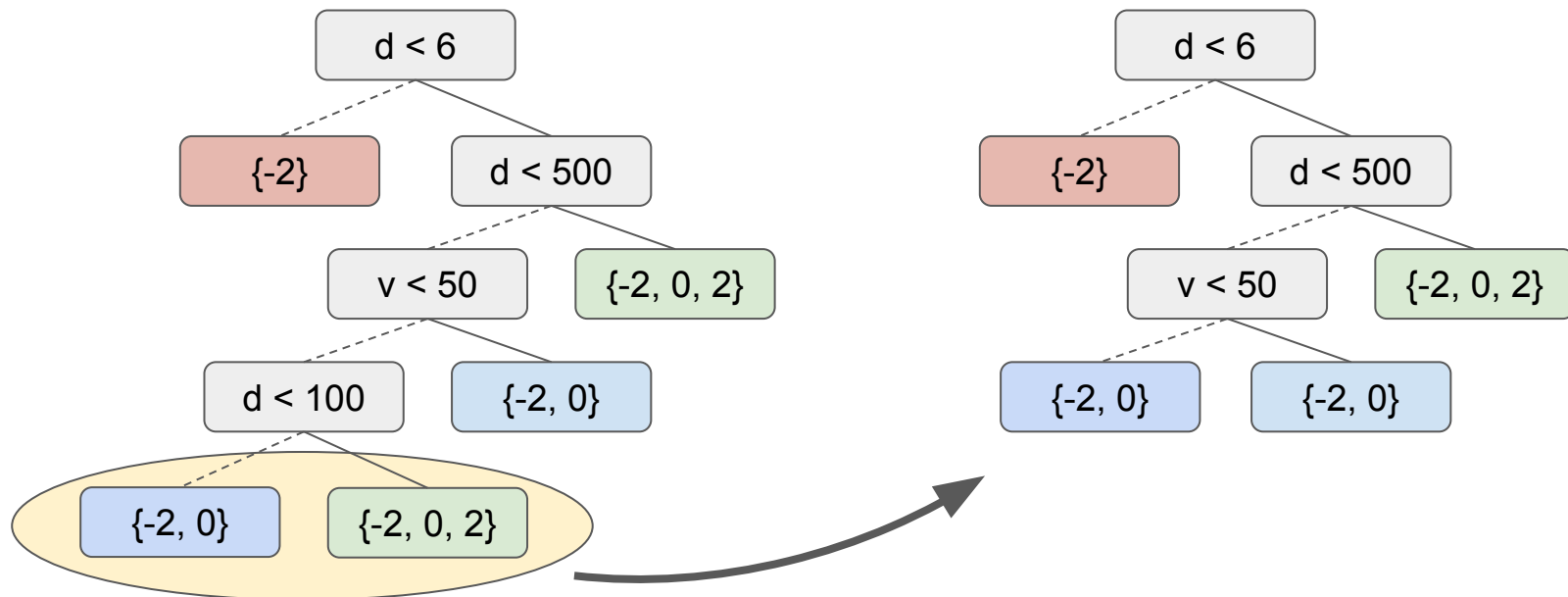{-2}

14

d < 500

9

{-2, 0, 2}

{-2}

{-2, 0, 2}

{0}

{-2, 0, 2}   {-2, 0}
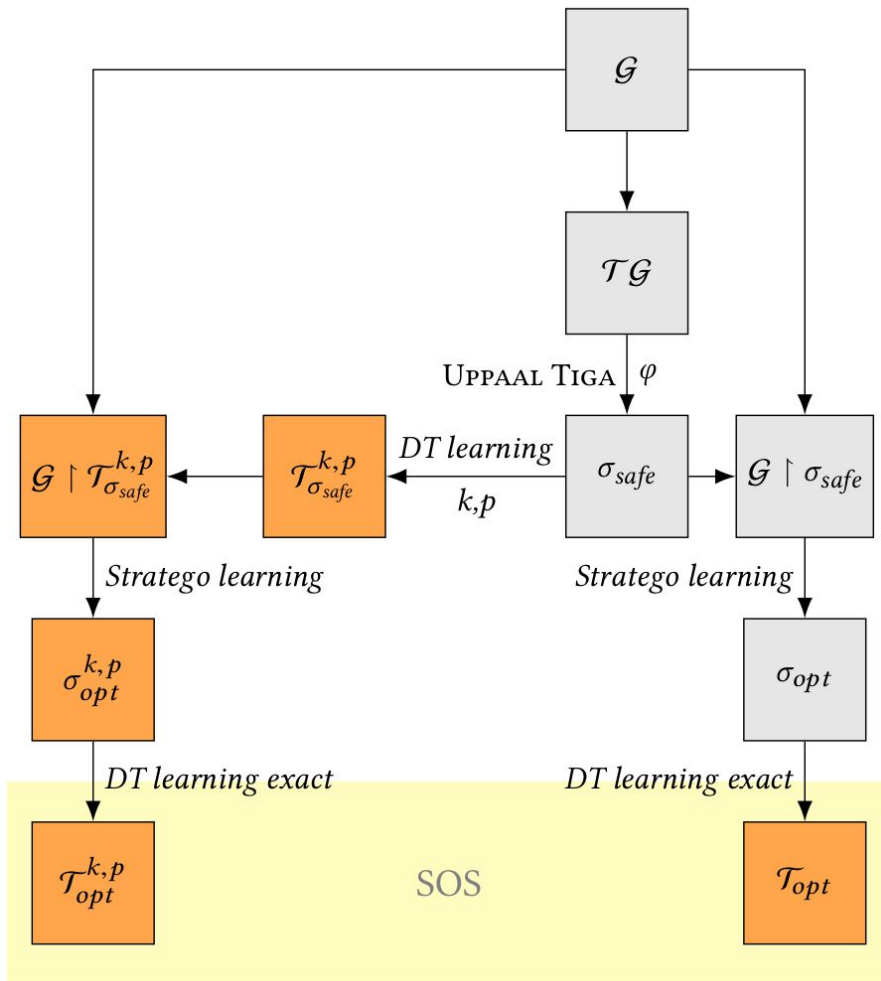
Consider splitting node only if # data points > k

# Safe Pruning

Merge leaves if intersection non-empty

Stratego+
Framework

# Experiments

**cruise**

    strategy guarantees safety only at integer points

**cruise-euler**

    enriched cruise

    strategy guarantees safety at all points

**tworooms-euler**

    automatic climate control

    interaction between two rooms, env. and heaters

# Experimental Results

| Model | State-action pairs | Controllable | BDD (reordered) | DT $T_{opt}$ |
|---|---|---|---|---|
| ~~cruise~~ | ~~1,790,034~~ | ~~308,216~~ | ~~5066~~ | ~~2,899~~ |
| ~~cruise-euler~~ | ~~5,931,154~~ | ~~304,752~~ | ~~4728~~ | ~~2,713~~ |
| cruise | 817,278 | 295,970 | 2,730 | 1,005 |
| cruise-euler | 1,140,756 | 414,899 | 2,667 | 1,025 |
| two-rooms | 1,924,708 | 509,715 | 20,214 | 487 |

| Min split size (k) | Rounds of pruning (p) | | | Min split size (k) | Rounds of pruning (p) | | |
|---|---|---|---|---|---|---|---|
| | 0 | 1 | 2 | | 0 | 1 | 2 |
| 2 | 2,713 | 1,725 | 1,267 | 2 | 2,627 | 3,618 | 4,240 |
| 10 | 2,705 | 1,733 | 1,249 | 10 | 2,696 | 3,596 | 4,210 |
| 20 | 2,667 | 1,733 | 1,131 | 20 | 2,778 | 3,625 | 14,039 |
| 30 | 2,657 | 1,695 | 993 | 30 | 2,778 | 3,589 | 14,108 |
| 40 | 2,627 | 1,669 | 1,015 | 40 | 2,778 | 3,600 | 14,096 |
| 50 | 2,557 | 1,695 | 1,003 | 50 | 2,825 | 3,614 | 14,037 |
| 60 | 2,635 | 1,489 | 963 | 60 | 2,905 | 3,673 | 14,074 |
| 70 | 2,613 | 1,441 | 955 | 70 | 2,898 | 3,714 | 14,095 |
| 80 | 2,519 | 1,537 | 915 | 80 | 2,907 | 3,717 | 14,092 |
| 90 | 2,455 | 1,323 | 923 | 90 | 3,006 | 3,741 | 14,077 |
| 100 | 1,929 | 1,023 | 877 | 100 | 3,030 | 14,061 | 14,292 |

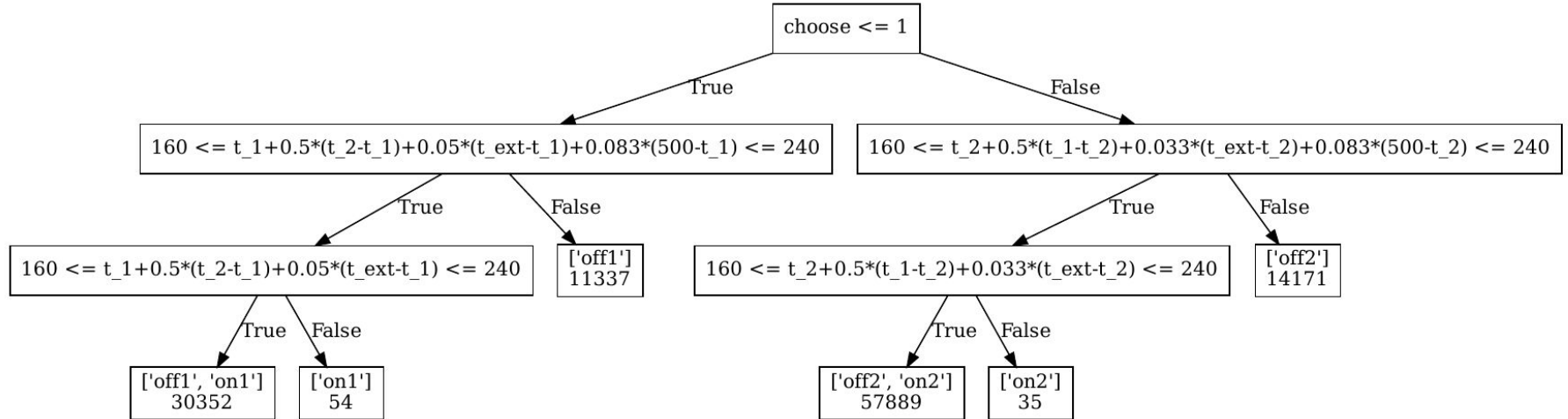Experimental Results: `cruise-euler` size-optimality tradeoff

# Handcrafted Strategy

```
1   // 0 - neutral, 1 - accelerate, 2 - decelerate
2   strategy(action, vE, d, vF, aF, aE) {
3       if (check(d, vE, vF, 2, aF) > 5) {
4           return action == 1;
5       } else if (check(d, vE, vF, 0, aF) > 5) {
6           return action == 0;
7       } else if (check(d, vE, vF, -2, aF) > 5) {
8           return action == 2;
9       } else {
10          return 0 == 1;
11      }
12  }
13
14  d_t(d, vE, vF, aE, aF, t) {
15      return d + 0.5*( aF - aE )*t*t + ( vF - vE )*t;
16  }
17
18  check(d, vE, vF, aE, aF) {
19      t1 = (vF + 10)/2;
20      if (t1 > 0.5) {
21          d1 = d_t(d, vE, vF, aE, -2, 1);
22          nvF = vF - 2;
23          nvE = vE + aE;
24      } else {
25          d1 = d_t(d, vE, vF, aE, 0, 1);
26          nvF = vF;
27          nvE = vE + aE;
28      }
29
30      if (t1 > 1) {
31          d2 = d_t(d1, nvE, nvF, -2, -2, t1 - 1);
32          nvE = nvE - 2*(t1 - 1);
33      } else {
34          d2 = d1;
35      }
36
37      t2 = (nvE + 10)/2;
38      if (t2 > 0) {
39          d3 = d_t(d2, nvE, -10, -2, 0, t2);
40      } else {
41          d3 = d2;
42      }
43
44      return d3;
45  }
```

# Strategy Preview: cruise

# Strategy Preview: tworooms

choose <= 1

True → 160 <= t_1+0.5*(t_2-t_1)+0.05*(t_ext-t_1)+0.083*(500-t_1) <= 240

False → 160 <= t_2+0.5*(t_1-t_2)+0.033*(t_ext-t_2)+0.083*(500-t_2) <= 240

From left branch:
- True → 160 <= t_1+0.5*(t_2-t_1)+0.05*(t_ext-t_1) <= 240
- False → ['off1'] 11337

From "160 <= t_1+0.5*(t_2-t_1)+0.05*(t_ext-t_1) <= 240":
- True → ['off1', 'on1'] 30352
- False → ['on1'] 54

From right branch:
- True → 160 <= t_2+0.5*(t_1-t_2)+0.033*(t_ext-t_2) <= 240
- False → ['off2'] 14171

From "160 <= t_2+0.5*(t_1-t_2)+0.033*(t_ext-t_2) <= 240":
- True → ['off2', 'on2'] 57889
- False → ['on2'] 35

# Concluding Remarks

**Problem**: Strategies from UPPAAL Stratego are large and incomprehensible

**Solution**: Stratego+ framework representing safe, small, optimal strategies as DTs

**Takeaways**

- BDDs are insufficient (uninterpretable, not that small either)
- Great prospects from decision trees

**Future work**

- Linear/algebraic predicates + domain knowledge

# Backup

# Experimental Results

| Model | State-action pairs | Controllable | BDD (median) | DT $T_{safe}$ | DT $T_{opt}$ |
|---|---|---|---|---|---|
| cruise | 817,278 | 295,970 | 2,730 | 1,017 | 1,005 |
| cruise-euler | 1,140,756 | 414,899 | 4,728 | 1,045 | 1,025 |
| two-rooms | 1,924,708 | 509,715 | 20,214 | 543 | 487 |