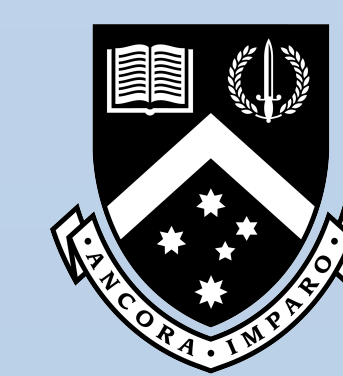# USING MaxSAT FOR EFFICIENT EXPLANATIONS OF TREE ENSEMBLES

Alexey Ignatiev[1], Yacine Izza[2], Peter J. Stuckey[1], Joao Marques-Silva[3]

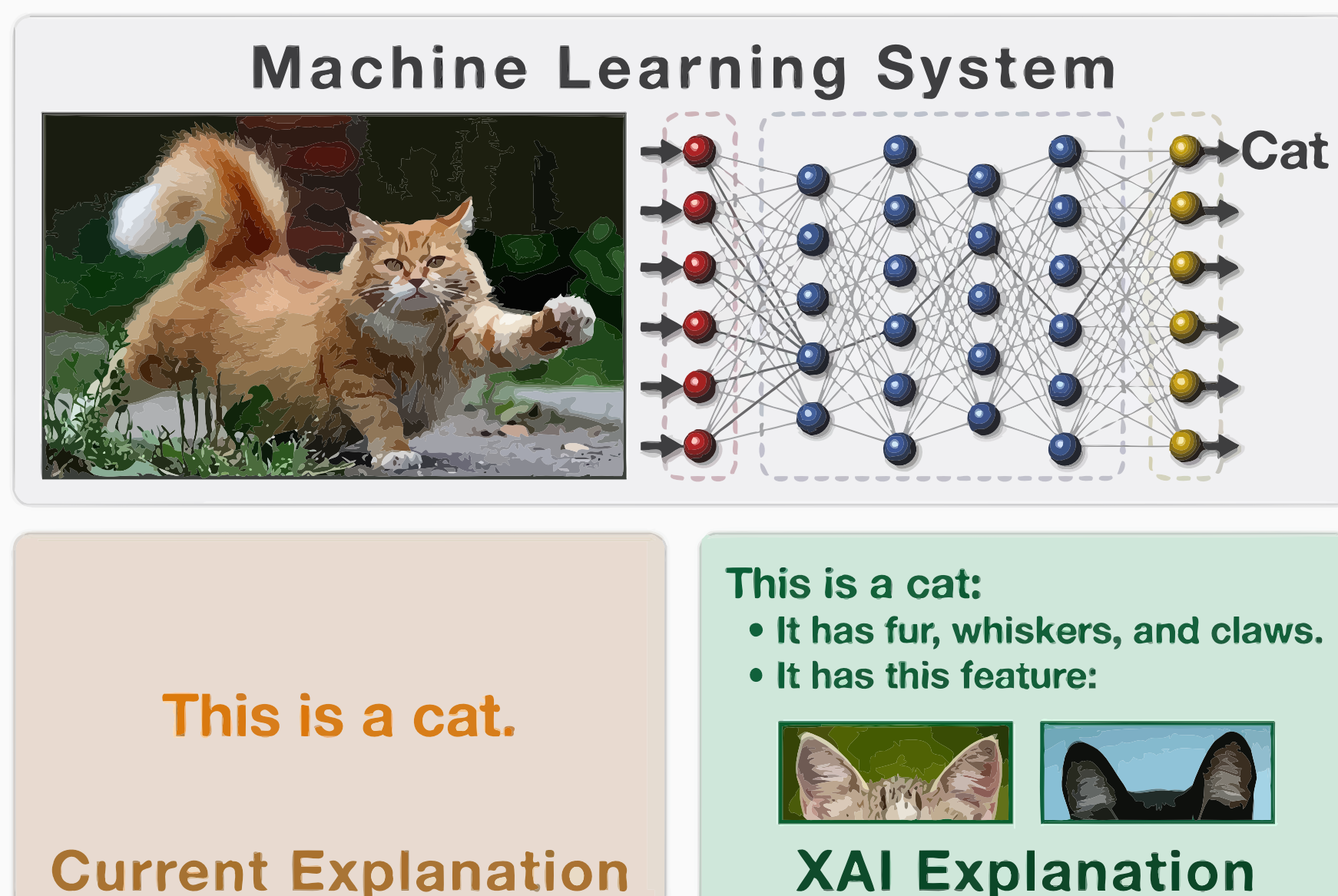[1]Monash University, Melbourne, Australia    [2]University of Toulouse, France    [3]IRIT, CNRS, Toulouse, France

MONASH University

ANITI · Artificial & Intelligence Toulouse Institute · Université Fédérale Toulouse Midi-Pyrénées

iRiT

## eXplainable AI



Machine Learning System

Cat

This is a cat.

Current Explanation

This is a cat:
• It has fur, whiskers, and claws.
• It has this feature:

XAI Explanation

## Why? Status Quo
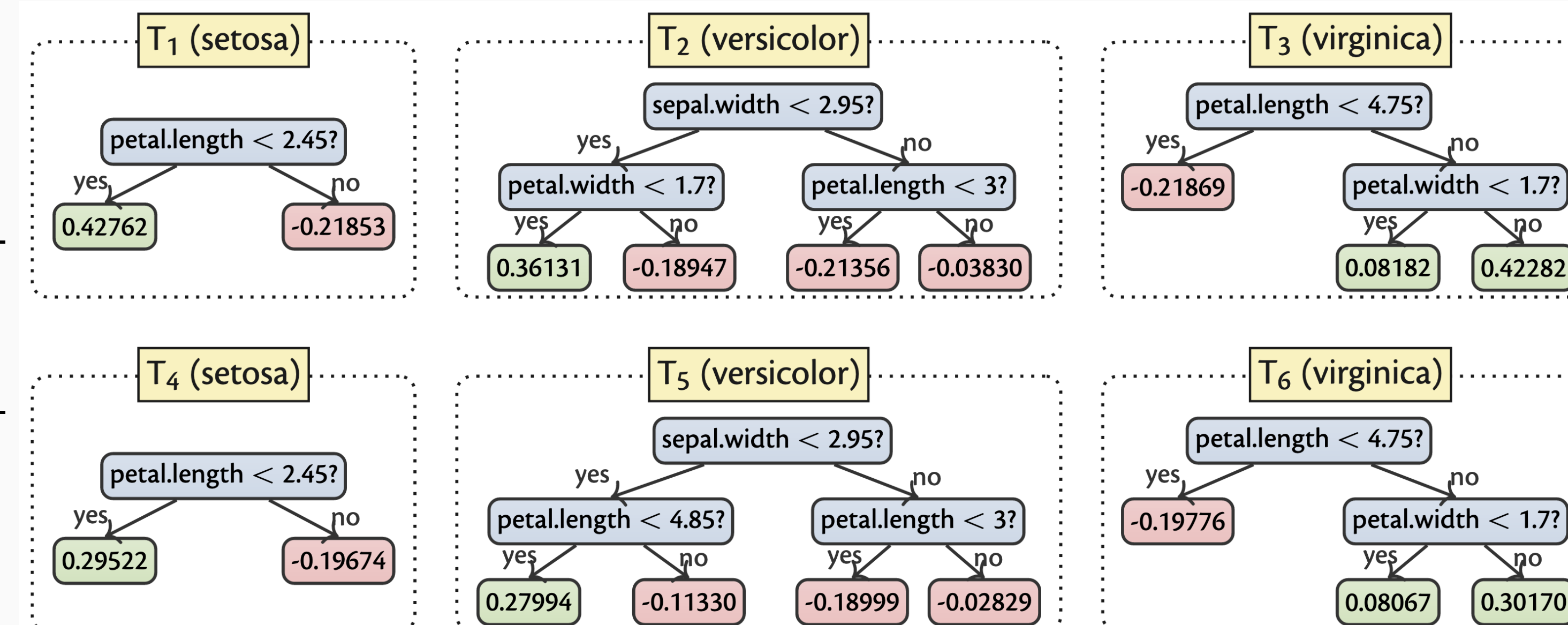
|  | A parrot | Machine learning algorithm |
|---|---|---|
| Learns random phrases | ✅ | ✅ |
| Doesn't understand s**t about what it learns | ✅ | ✅ |
| Occasionally speaks nonsense | ✅ | ✅ |

**classifier** $\tau : \mathbb{F} \to \mathcal{K}$,
**instance** $\mathbf{v}$ s.t. $\tau(\mathbf{v}) = c$

**abductive explanation** $\mathcal{X}$:
$\forall(\mathbf{x} \in \mathbb{F}) . \bigwedge_{j \in \mathcal{X}}(x_j = v_j) \to (\tau(\mathbf{x}) = c)$

• $w(\mathbf{x}, c) = \sum_{j \in \{0,...,n-1\}} \mathcal{T}_{Kj+c}(\mathbf{x}),\ c \in [K]$
• $\tau(\mathbf{x}) = \arg\max_{c \in [K]} w(\mathbf{x}, c)$

## Abductive Explanations and Boosted Trees



$T_1$ (setosa), $T_2$ (versicolor), $T_3$ (virginica), $T_4$ (setosa), $T_5$ (versicolor), $T_6$ (virginica)

## Idea and Contributions

**SMT**
• reach logic, can handle *linear constraints*
• directly reason about:
$\mathcal{H} \wedge \left(\sum_{\iota} \geq \sum_i\right)$
• high decimal point precision

**MaxSAT**
**maximize** $\sum_{\iota} - \sum_i$
**subject to** $\mathcal{H}$

• pure propositional logic
  – prediction as objective function
  – weighted soft clauses keep precision
  – core-guided MaxSAT
• incremental MaxSAT calls
  – *MiniSat-like assumptions interface!*
  – unsatisfiable core *reuse*
• novel weight stratification
• early termination

## Encoding BT Operation

$\mathcal{F}$ — set of features    $\forall_{j \in \mathcal{F}}\ D_j$ — domain of feature $j$    $\mathfrak{E}$ — tree ensemble

$\mathcal{H} = \bigwedge_{j \in \mathcal{F}} \mathcal{H}_{D_j} \wedge \bigwedge_{T_i \in \mathfrak{E}} \mathcal{H}_{T_i}$

$\mathcal{H}_{D_j}$ encodes feature domain $D_j$
• feature **threshold values** $s_{j,k}$ from $\mathfrak{E}$
• variable $o_{j,k} = 1$ iff $x_j < s_{j,k}$
• variable $l_{j,k} = 1$ iff $x_j \in [s_{j,k-1}, s_{j,k})$
• **order encoding** of $D_j$ with $l_{j,k}$ and $o_{j,k}$
• instance is expressed by $l_{j,k}$-variables

$\mathcal{H}_{T_i}$ encodes paths of $T_i$
• $\mathcal{P}_i$ is the set of paths of $T_i$
• given a $P_r \in \mathcal{P}_i$, $t_r$ denotes its leaf
• encode each path $P_r \in \mathcal{P}_i$:
$$\left(\bigwedge_{(x_j < s_{j,k}) \in P_r} o_{j,k} \wedge \bigwedge_{(x_j \geq s_{j,k}) \in P_r} \neg o_{j,k}\right) \leftrightarrow t_r$$
• $\sum_{P_r \in \mathcal{P}_i} t_r = 1$

## "Optimizing" Model Predictions

$\forall(\mathbf{x} \in \mathbb{F}). \bigwedge_{j \in \mathcal{X}}(x_j = v_j) \to (\tau(\mathbf{x}) = c_i)$

AXp condition for $\mathcal{X}$

$\mathcal{X}$ is not AXp if $\exists_{\mathbf{x} \in \mathbb{F}, c_i \neq c_i}. w_i(\mathbf{x}) \geq w_i(\mathbf{x})$

weighted sums $\sum_{\iota}$ and $\sum_i$

even to represent in propositional logic!    hard to reason about

$\mathcal{X}$ is not AXp if $\mathcal{S}_{\iota,\iota}^* \geq 0$

**maximize** $\mathcal{S}_{\iota,\iota} = \sum_{\iota} - \sum_i$
**subject to** $\mathcal{H} \wedge \bigwedge_{j \in \mathcal{X}} \mathbb{1}[x_j = v_j]$

weighted soft clauses + propositional variables only

## Computing an AXp

```
Function EntCheck(⟨H, S⟩, v, c_i, X)
  Input:  H: Hard clauses,   S: Objective functions,
          v: Input instance,  c_i: Prediction, i.e. τ(v) = c_i,
          X: Candidate explanation
  Output:  true or false
1   foreach S_{i,t} ∈ S:                  # all relevant objective functions for c_i
2     μ ← MaxSAT(H ∧ ⋀_{j∈X} l[[x_j = v_j]], S_{i,t})
3     if ObjValue(μ) ≥ 0:                 # non-negative objective?
4       return false                       # misclassification reached
5   return true                            # X indeed entails prediction c_i

Function ExtractAXp(𝔈, v, c_i)
  Input:  𝔈: TE computing τ(x),  v: Input instance,
          c_i: Prediction, i.e. τ(v) = c_i
  Output:  X: abductive explanation
1   ⟨H, S⟩ ← Encode(𝔈)                    # MaxSAT encoding s.t. S ≜ ∪S_{i,t}
2   X ← F                                  # X is over-approximation
3   foreach j ∈ X:
4     if EntCheck(⟨H, S⟩, v, c_i, X \ {j}): # j unneeded?
5       X ← X \ {j}                         # If so, drop it
6   return X                               # X is AXp
```

## Incremental Core-Guided MaxSAT Solver

```
Function MaxSAT(φ, A)
  Input:  φ: Partial CNF formula,
          A: Set of assumption literals
  Output:  μ: MaxSAT model
1   cost ← 0                              # initially, cost is 0
2   C ← ValidCores(φ, A)                  # get valid unsatisfiable cores
3   foreach κ ∈ C:                        # iterate over known cores κ
4     cost ← cost + CoreWt(κ)             # add its weight to cost
5     φ ← Process(φ, κ)                   # process κ and update φ
6   while SAT(φ, A) = false:              # iterate until φ gets satisfiable
7     κ ← GetCore(φ)                      # new unsatisfiable core
8     cost ← cost + CoreWt(κ)             # add its weight to cost
9     φ ← Process(φ, κ)                   # process κ and update φ
10    Record(φ, A, κ)                     # record κ for the future
11  return GetModel(φ)                    # φ is now satisfiable
```

## Experimental Results



**all three competitors**



**MaxSAT vs. SMT**



**MaxSAT vs. Anchor**

up to 2 orders of magnitude performance improvement

more robust than SMT and Anchor