# From Contrastive to Abductive Explanations and Back Again

Alexey Ignatiev[1], Nina Narodytska[2], Nicholas Asher[3], and Joao Marques-Silva[3]

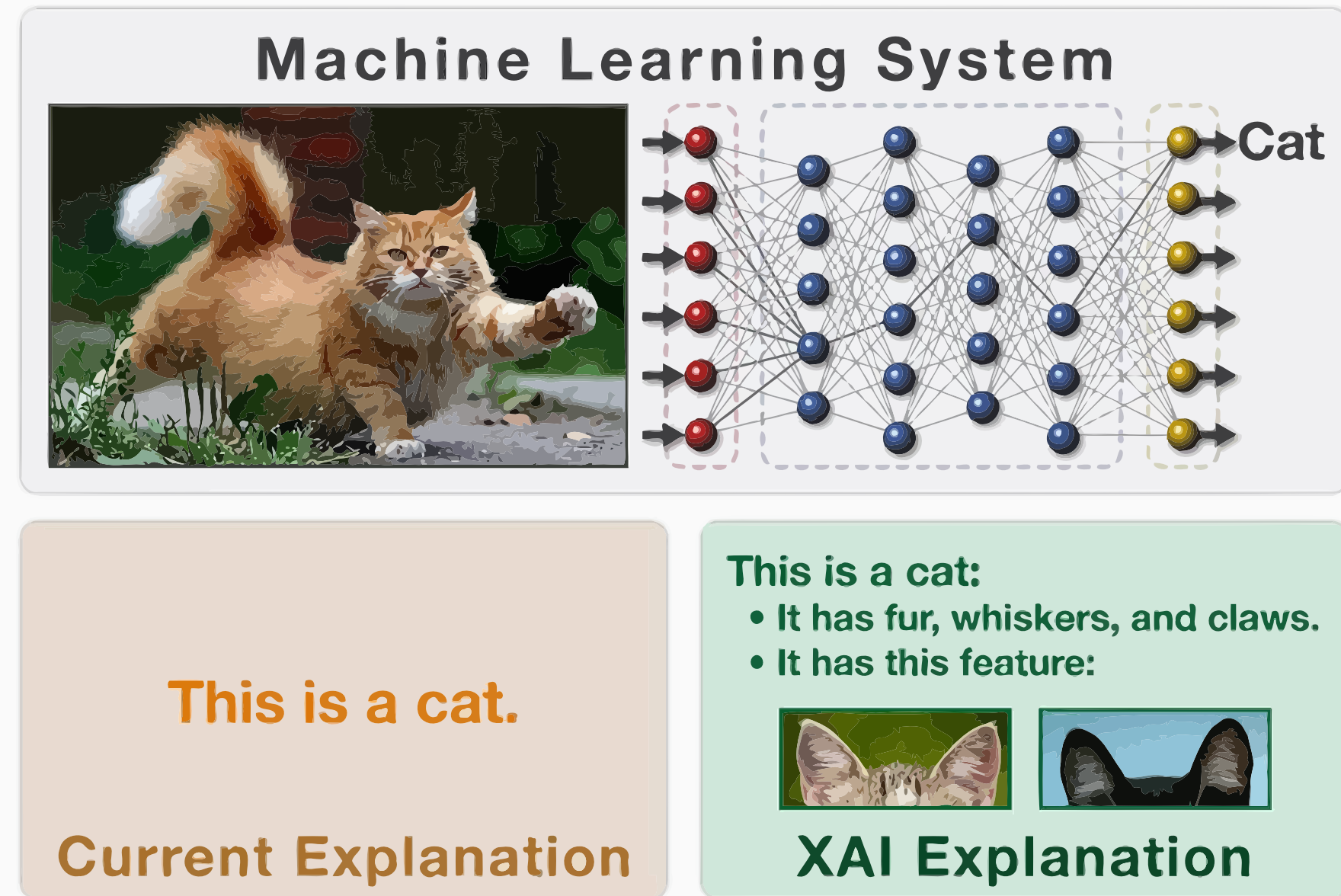[1]Monash University, Melbourne, Australia   [2]VMware Research, CA, USA   [3]IRIT, CNRS, Toulouse, France
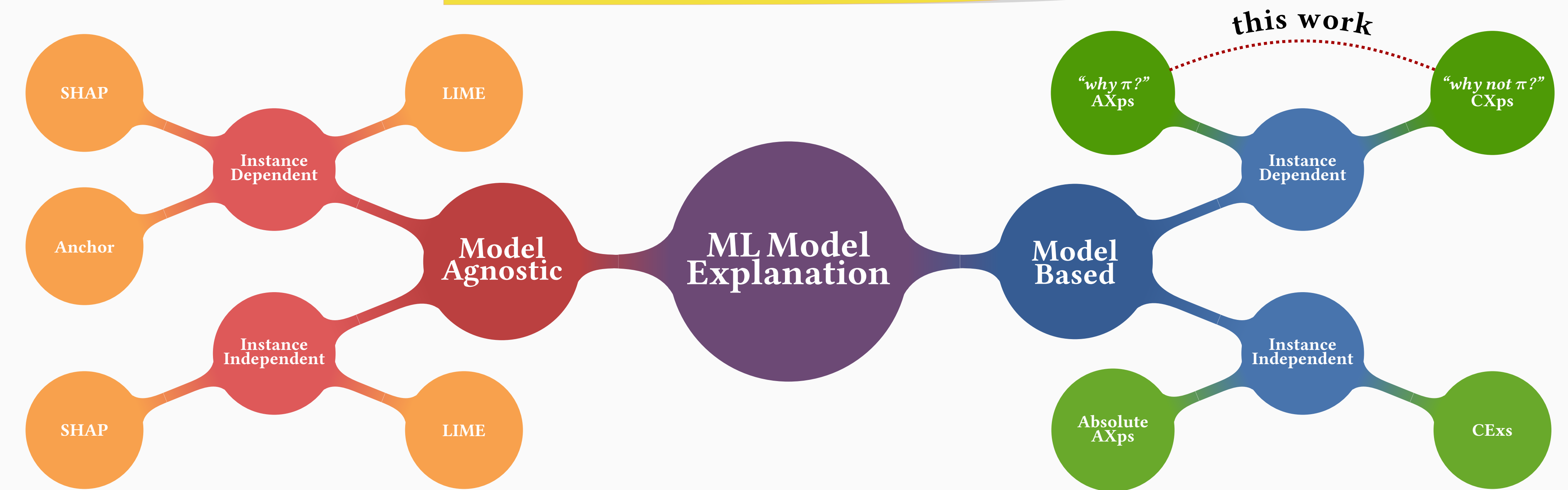
MONASH University    vmware®    iRIT

## 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 | ☑ | ☑ |

## Taxonomy of ML Model Explanations



*this work*

SHAP — LIME — Instance Dependent — Model Agnostic — ML Model Explanation — Model Based — Instance Dependent — "why π?" AXps / "why not π?" CXps

Anchor — Instance Independent

SHAP — LIME — Instance Independent

Absolute AXps — Instance Independent — CExs

## Formal Explanations

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

abductive explanation $\mathcal{X}$

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

contrastive explanation $\mathcal{Y}$

$$\exists (\mathbf{x} \in \mathbb{F}) . \bigwedge_{j \notin \mathcal{Y}} (x_j = v_j) \wedge (\tau(\mathbf{x}) \neq c)$$

## Explanation Examples

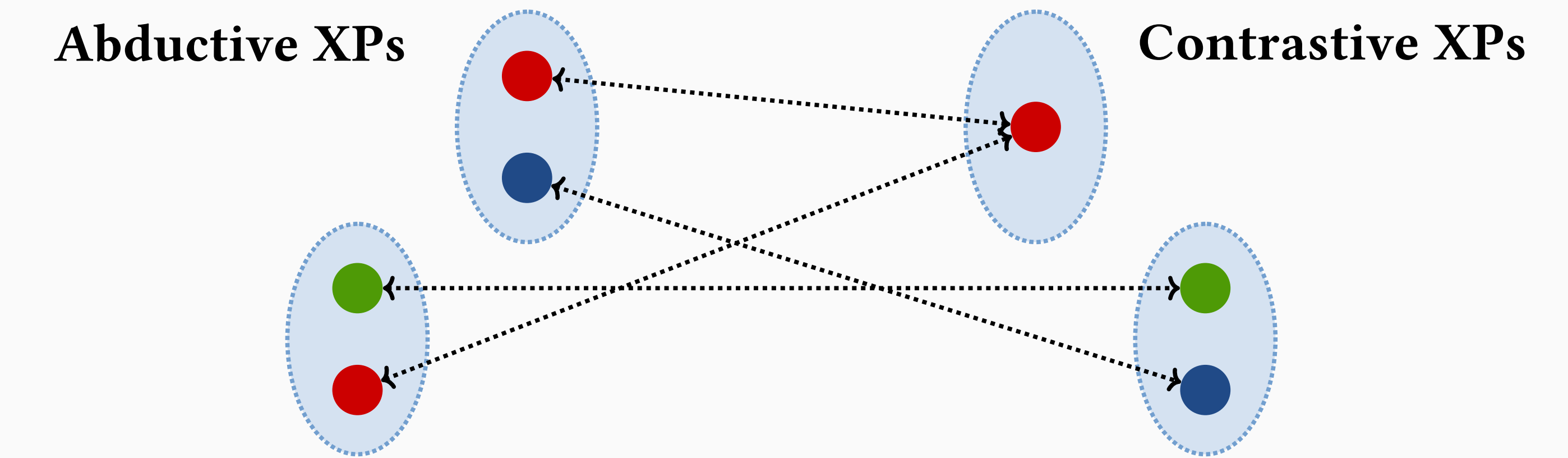$$\mathbb{F} = \{0, 1, 2\}^5 \qquad \mathcal{K} = \{\ominus, \oplus\}$$

| $R_0$: | IF | $x_1 = 1 \wedge x_2 = 1$ | THEN $\ominus$ |
|---|---|---|---|
| $R_1$: | ELSE IF | $x_3 \neq 1$ | THEN $\oplus$ |
| $R_{\text{DEF}}$: | ELSE | | THEN $\ominus$ |

**observe** $\tau(1, 1, 1, 1, 1) = \ominus$

AXps $\mathbb{X} = \{\{1, 2\}, \{3\}\}$
CXps $\mathbb{Y} = \{\{1, 3\}, \{2, 3\}\}$

## Minimal Hitting Set Duality

**Abductive XPs**          **Contrastive XPs**



**AXps are minimal hitting sets of CXps, and vice versa.**

## Enumerating CXps

```
Function CXpEnum(τ, v, c)
  Input:  τ: ML model, v: Input
          instance, c = τ(v): Prediction
1   I ← ∅                    // Block CXps
2   while true:
3     μ ← ExtractCXp(τ, v, c, I)
4     if μ = ∅: break
5     ReportCXp(μ)
      I ← I ∪ ⋁_{j∈μ} (x_j = v_j)

Function ExtractCXp(τ, v, c, I)
  Input:  τ: classifier, v: Input instance,
          c = τ(v): Prediction, I: Blocked
          CXps
  Output: S: Minimal set
1   S ← [|v|]
2   foreach j ∈ S:
3     if SAT(⋀_{i∈S\{j}}(x_i = v_i) ∧ I ∧ τ(x) ≠ c):
4       S ← S \ {j}
5   return S              // S is CXp
```

## Enumerating AXps and CXps

```
Function XpEnum(τ, v, c)
  Input:  τ: ML model, v: Input instance, c = τ(v):
          Prediction
1   K = (N, P) ← (∅, ∅)          // Block AXps & CXps
2   while true:
3     (st_λ, λ) ← FindMHS(P, N)   // MHS of P s.t. N
4     if ¬st_λ: break
5     st_c' ← SAT(⋀_{j∈λ}(x_j = v_j) ∧ τ(x) ≠ c)
6     if ¬st_c':                  // entailment holds
7       ReportAXp(λ)
8       N ← N ∪ ⋁_{j∈λ}(x_j ≠ v_j)
9     else:
10      μ ← ExtractCXp(τ, v, c, P)
11      ReportCXp(μ)
12      P ← P ∪ ⋁_{j∈μ}(x_j = v_j)
```

## Explanation Enumeration Results

| | | Dataset | | | | |
|---|---|---|---|---|---|---|
| | Adult | Lending | Recidivism | Compas | German | Spambase |
| # of instances | 5579.0 | 4414.0 | 3696.0 | 778.0 | 1000.0 | 2344.0 |
| total time (sec.) | 7666.9 | 443.8 | 3688.0 | 78.4 | 16 943.2 | 6859.2 |
| minimal time (sec.) | 0.1 | 0.0 | 0.1 | 0.0 | 0.2 | 0.1 |
| average time (sec.) | 1.4 | 0.1 | 1.0 | 0.1 | 16.9 | 2.9 |
| maximal time (sec.) | 13.1 | 0.8 | 8.9 | 0.5 | 193.0 | 23.1 |
| total oracle calls | 492 990.0 | 69 653.0 | 581 716.0 | 21 227.0 | 748 164.0 | 176 354.0 |
| minimal oracle calls | 14.0 | 11.0 | 17.0 | 13.0 | 23.0 | 12.0 |
| average oracle calls | 88.4 | 15.8 | 157.4 | 27.3 | 748.2 | 75.2 |
| maximal oracle calls | 581.0 | 73.0 | 1426.0 | 134.0 | 7829.0 | 353.0 |
| total # of AXps | 52 137.0 | 8105.0 | 60 688.0 | 1931.0 | 59 222.0 | 18 876.0 |
| average # of AXps | 9.4 | 1.8 | 16.4 | 2.5 | 59.2 | 8.1 |
| average AXp size | 5.3 | 1.9 | 6.4 | 3.8 | 7.5 | 4.6 |
| total # of CXps | 66 219.0 | 8663.0 | 77 784.0 | 3558.0 | 66 781.0 | 24 774.0 |
| average # of CXps | 11.9 | 2.0 | 21.1 | 4.6 | 66.8 | 10.6 |
| average CXp size | 2.4 | 1.4 | 2.6 | 1.5 | 3.6 | 2.3 |

## Debugging SHAP



Real 6 | XGBoost | SHAP | CXp¹ | CXp² | CXp¹⁻³

Fake 6 | XGBoost | SHAP | CXp³ | CXp⁴ | CXp³⁻⁵

The **"real vs fake"** images. The first row shows results for the *real image 6*; the second – results for the *fake image 6*. The first column shows examples of inputs; the second – heatmaps of XGBoost's important features; the third – heatmaps of SHAP's explanation. Last three columns show heatmaps of CXp of different cardinality. The brighter pixels are more influential features.