DEEL

# FROM END 2018…



30 — Million of total Budget
27 — Partners
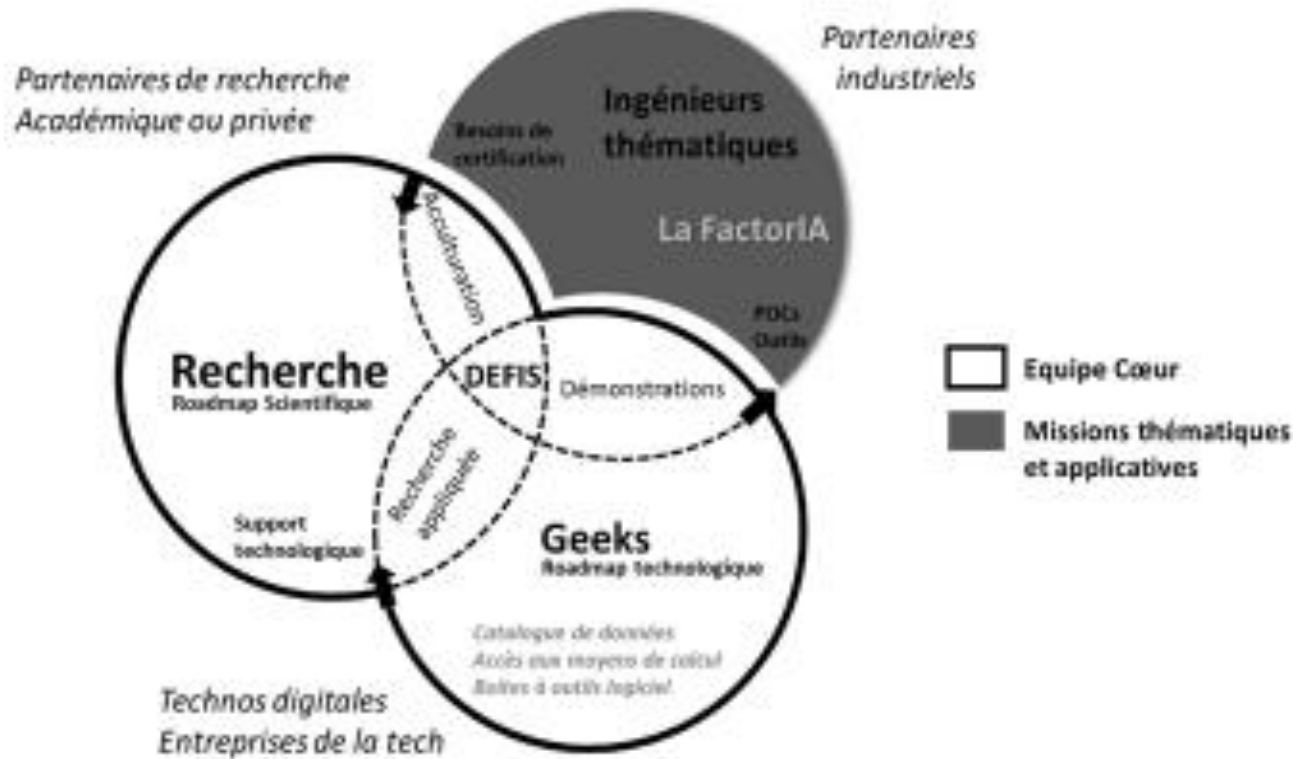50 — Researchers, geeks & thematic engineers
2 — Countries
25 — PhDs

# …TO END 2023

# ORGANISATION FR

Principe de colocalisation des équipes



Antenne à Montréal

# EXPLAINABILITY CHALLENGE

# EXPLAINABILITY CHALLENGE

- Kickoff : 2019 October 31th
- Team: 7 persons + 1 Phd + 2 researchers

# OUTPUT

- State of the art document
  - Audiences (Who ?)
  - Explainability type of output (What ?)
  - Problematics (Why ?)
  - Design process (When ?)
  - Industrial Uses-cases
  - Toolboxes
  - **Mapping between the technics and the audience, problematics and industrial usecases**
  - State of the art description
- Current Results:
  - Explainability toolbox
    - Local
    - Global
    - Metrics
  - Notebooks to evaluate technics

# WHY DO WE NEED EXPLANATIONS

# Why do we need Explanations

Build trust in the model prediction [3][4]

Elucidate important aspects of learned models [4]

Help satisfy regulatory requirements and Certification process[1]

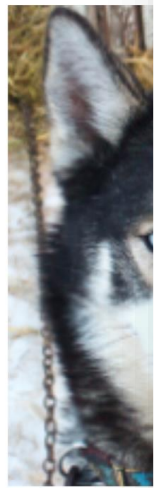Reveal bias or other unintended effects learned by a model [3]

[1] Bryce Goodman & al. European union regulations on algorithmic decision-making and a"right to explanation".
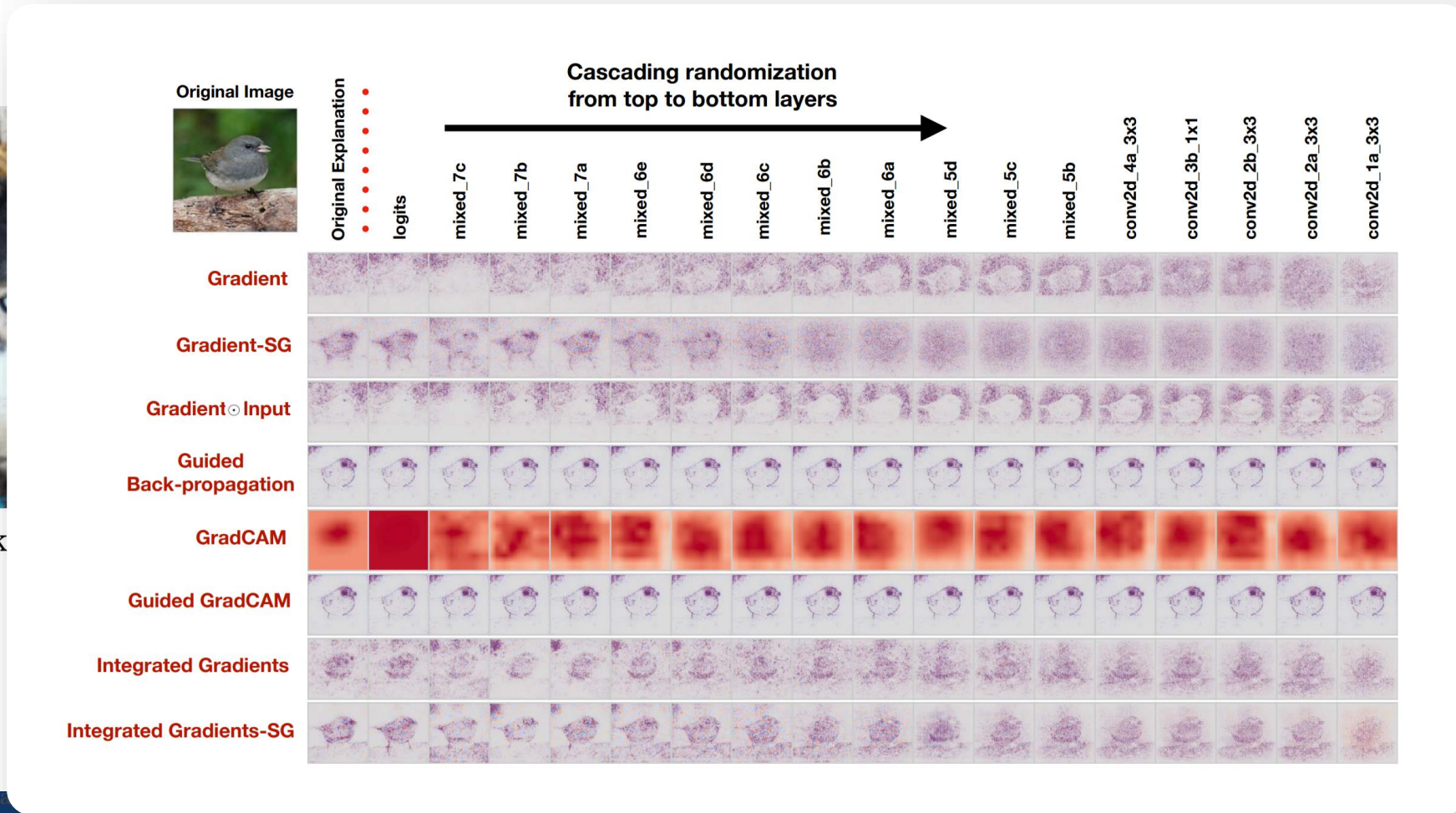[2] Finale Doshi-Velez & al. Accountability of ai under the law: The role of explanation
[3] Gabriel Cadamuro & al. Debugging machine learning models
[4] Alfredo Vellido & al. Making machine learning models interpretable

# "What is a good explanation?"



(a) Husk...

Sa...
Explaining the Predictions of Any Classifier (2016)
Evaluating the visualization of what a Deep Neural Network has learned (2016)

Confirmation bias.

**Just because it makes sense to humans doesn't mean it reflects the evidence for prediction.**

# STATE OF THE ART

# STATE OF THE ART OVERVIEW (1/2)

- Global explanations
    - Transparency models
    - Features relevance explanations
    - Explanation by simplification
    - Internal analysis
    - Explanation by examples
    - Natively explainable models
        - Models providing an explanation as output
        - Building interpretable features
        - Attention models
        - Unsupervised learning for representation disentanglement

# STATE OF THE ART OVERVIEW (2/2)

- Global explanations:
  - Causality
  - Formal methods

- Local explanation

- Validation:
  - Metrics
  - Explainability Robustness
  - Link between Robustness and Explicability

# WORKING AXES

# 2 MAIN AXES

- Research thematics:
  - Goal: Develop research axes which are important for Deel project and which are not much investigated in the research community

- Deel Explainability "Library": Evaluation of existing technologies on our industrial usecases
  - Goal: Create software suite and Jupyter notebook tutorials
    - Tutorials are given to explain how Explainability techniques shall be used to analyse different industrial uses cases
    - The techniques could be implemented in a DEEL library or relying on existing external toolboxes

# RESEARCH THEMATICS

- Outcomes : Articles & source code
- **Metrics / Explainability Robustness:**
  - 2 core team members
  - 1 researcher
  - 1 PhD student
- **Formal methods**
  - 2 core team members
  - Link to ANITI project
- Backlog:
  - Causality
  - Link between Robutness and Explainability
  - Building interpretable features & Attention models

Library for Global Explaination:
www.gems-ai.com

# DEEL EXPLAINABILITY "LIBRARY"

- Outcome:
  - Deel Explainability Library (source code)
  - Tutorials (Jupyter notebook)
  - Feed back on industrial usecases

| | | |
|---|---|---|
| **1** | Internal Analysis | |
| **2** | Building features/ attentions / Unsupervised learning for representation disentanglement | |
| **3** | Formal methods | |
| **4** | Inputs local Importance | |
| **5** | Causality | |
| **6** | Link between explicability / Robustness | |
| **7** | Metrics | |

Internal model analysis



First evaluation done

3 core team members

2 core team members

VAE Evaluation on Deel Dataset

Amount of pins



90 degree rotation

# REPRESENTATIVITY AND CONSISTENCY MEASURES FOR DEEP NEURAL NETWORK EXPLANATIONS

# Properties of explainability

no metrics associated

**Fidelity**            Does my explanation reflect the behavior of my model?

**Representativity**    How many phenomena my explanation cover?

**Comprehensibility**   Is my explanation unambiguous and simple?

**Consistency**         The degree to which similar explanations are generated from different models trained on the same task.

**Stability**           Does my explanation remain the same under semantically invariant transformation?

**Novelty**             Does my explanation reflect the fact that explained instance is from a new region, not contained or well represented in the training set?

Explanation in Artificial Intelligence: Insights from the social sciences **(2019)**
Explanation Methods in Deep Learning: Users, Values, Concerns and Challenges **(2018)**
Evaluating Explanation Without Ground Truth in Interpretable Machine Learning **(2019)**
Perturbation-Based Explanations of Prediction Models **(2018)**

**Consistency**

An explanation leading to predict **y** and **¬y** is inconsistent.

**Representativity**

A model should not base an explanation on a single sample.

Split the dataset



$$\forall (x, y) \in \mathcal{D}, \ \forall \mathcal{D}_i : x \in \mathcal{D}_i, \ \forall \mathcal{D}_j : x \notin \mathcal{D}_j$$

$$\mathcal{T} = \{d(\phi_x^{\mathcal{D}_i}, \phi_x^{\mathcal{D}_j}) \mid f_{\mathcal{D}_i}(x) = y \wedge f_{\mathcal{D}_j}(x) = y\}$$

$$\mathcal{Z} = \{d(\phi_x^{\mathcal{D}_i}, \phi_x^{\mathcal{D}_j}) \mid f_{\mathcal{D}_i}(x) = y \oplus f_{\mathcal{D}_j}(x) = y\}$$

# RELATIVE CONSISTENCY

$$ReCo = \max_\gamma \; TPR(\gamma) + TNR(\gamma) - 1$$



# MEAN GENERALIZABILITY

$$MeGe = \frac{1}{1+\mathbb{E}[\mathcal{T}]}$$

ReCo scores for normally trained ResNet-18

| Model - Dataset | IG | SG | SA | GI | GC |
|---|---|---|---|---|---|
| Cifar10 | 0.107 | 0.154 | 0.151 | 0.088 | **0.637** |
| Cifar100 | 0.018 | 0.132 | 0.131 | 0.004 | **0.800** |
| EuroSAT | 0.309 | 0.182 | 0.177 | 0.241 | **0.591** |
| FashionMNIST | 0.369 | 0.125 | 0.1 | 0.369 | **0.517** |

**Fashion MNIST**

| | IG | SG | SA | GI | GC | Shap |
|---|---|---|---|---|---|---|
| **Lipschitz** | **0.598** | **0.898** | **0.81** | **0.5** | **0.668** | 0.38 |
| ResNet-18 | 0.107 | 0.154 | 0.151 | 0.088 | 0.637 | **0.387** |

*Table 2. ReCo* score obtained on Cifar10. Higher is better.

Preprint:

08/04/2021

# MERCI POUR VOTRE ATTENTION