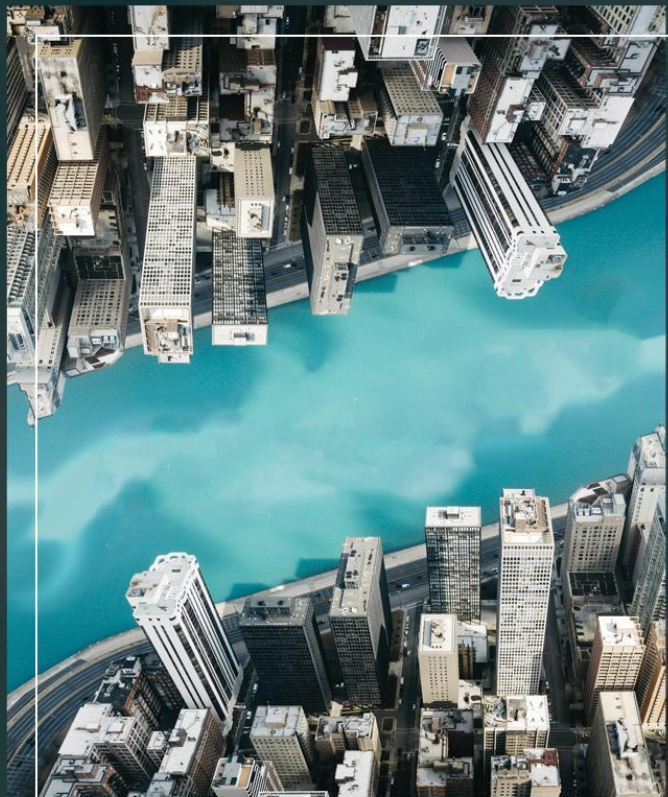# Concept-based Explainability:
## Challenges & Applications to Fraud Detection

Vladimir Balayan, Catarina Belém, Pedro Saleiro, Ludwig Krippahl, Pedro Bizarro

Deep Learning Sessions Lisbon - Meetup
2021, July 21st

# Problem & Motivation

# Problem

- The field of eXplainable AI (XAI) aims to tackle the lack of interpretability in ML.

# Problem

- The field of eXplainable AI (XAI) aims to tackle the lack of interpretability in ML.

- State-of-the-art methods in explainable AI (XAI) either:

| **Motivation** | Related Work | Solution | Experiment | Conclusion |

# Problem

- The field of eXplainable AI (XAI) aims to tackle the lack of interpretability in ML.

- State-of-the-art methods in explainable AI (XAI) either:

  1. produce low-level feature attributions explanations that are not suited for non-ML experts (e.g. fraud analyst).

# Problem

- The field of eXplainable AI (XAI) aims to tackle the lack of interpretability in ML.

- State-of-the-art methods in explainable AI (XAI) either:

  1. produce low-level feature attributions explanations that are not suited for non-ML experts (e.g. fraud analyst).

  Or

  2. produce concept-based explanations that do not work for tabular data.

| Motivation | Related Work | Solution | Experiment | Conclusion |

# Domain expert reasoning example

# Domain expert reasoning example

There are **multiple shipping addresses** in the last hour, so it seems **reshipping...** It's fraud!

Fraud Analyst

The ***ideal*** human-interpretable explanation for domain experts provide the ***high-level*** insights about the models' predictions.

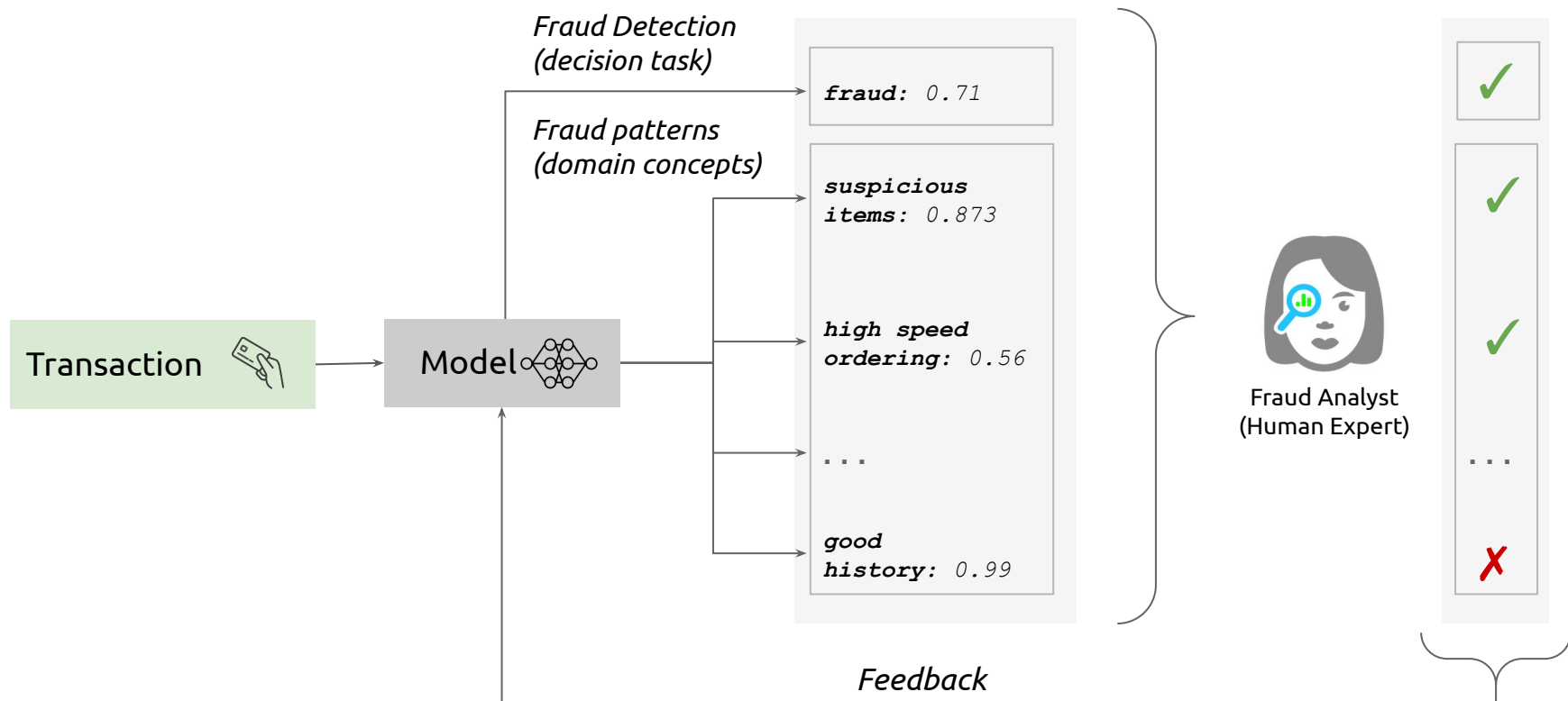| Motivation | Related Work | Solution | Experiment | Conclusion |

# Goals

- Develop a self-explainable neural network that jointly learn a predictive task and also associated domain knowledge explanations.

| **Motivation** | Related Work | Solution | Experiment | Conclusion |

# Goals

- Develop a self-explainable neural network that jointly learn a predictive task and also associated domain knowledge explanations.

- Develop a taxonomy of fraud concepts to be used as explanations.

| Motivation | Related Work | Solution | Experiment | Conclusion |

# Goals

- Develop a self-explainable neural network that jointly learn a predictive task and also associated domain knowledge explanations.

- Develop a taxonomy of fraud concepts to be used as explanations.

- Leverage the human-in-the-loop feedback to continuously improving both predictive accuracy and explainability.

| **Motivation** | Related Work | Solution | Experiment | Conclusion |

# Goals

- Develop a self-explainable neural network that jointly learn a predictive task and also associated domain knowledge explanations.

- Develop a taxonomy of fraud concepts to be used as explanations.

- Leverage the human-in-the-loop feedback to continuously improving both predictive accuracy and explainability.

- Create a semantic mapping bootstrapping strategy for automatically labeling concept-based explanations dataset.

| Motivation | Related Work | Solution | Experiment | Conclusion |

# Proposed solution in a real world fraud detection setting

feedzai



Transaction → Model

*Fraud Detection (decision task)* → **fraud:** *0.71*

*Fraud patterns (domain concepts)*

**suspicious items:** *0.873*

**high speed ordering:** *0.56*

`...`

**good history:** *0.99*

✓

✓

✓

...

✗

Fraud Analyst
(Human Expert)

*Feedback*

# Background & Related Work

# XAI Personas - Different personas, different XAI needs...

**Fields of expertise:**

- Domain knowledge
- No ML knowledge

- Data Science + ML

- Know what they want
- No domain nor ML knowledge

- Regulations/law
- Limited domain & ML knowledge

### Human-in-the-loop

### Data Scientist

### Decision Subject

### Regulator

**Goals:**

- Efficiency
- Better & faster decisions

- Efficiency
- Iterate and debug models

- Improve outcomes
- Reduce friction while feeling safe

- Audit and Assess
- if the system is compliant.

# XAI Personas - Different personas, different XAI needs...

feedzai

**Fields of expertise:**

- Domain knowledge
- No ML knowledge

- Data Science + ML

- Know what they want
- No domain nor ML knowledge

- Regulations/law
- Limited domain & ML knowledge

### Human-in-the-loop

### Data Scientist

### Decision Subject
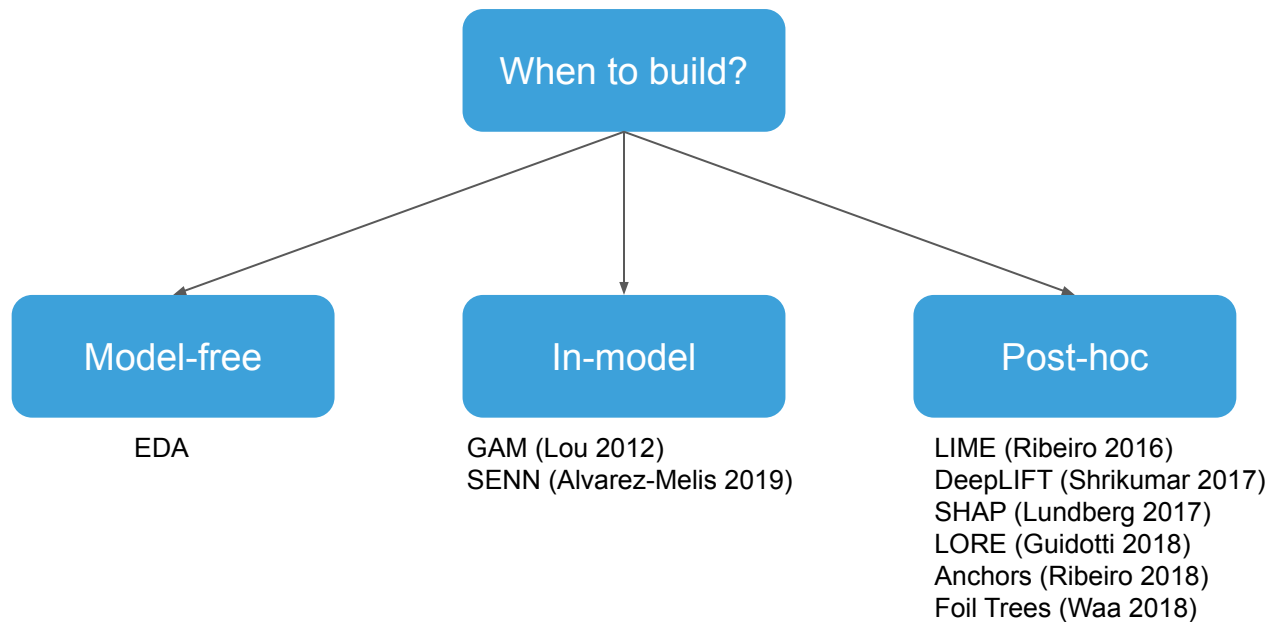
### Regulator

**Goals:**

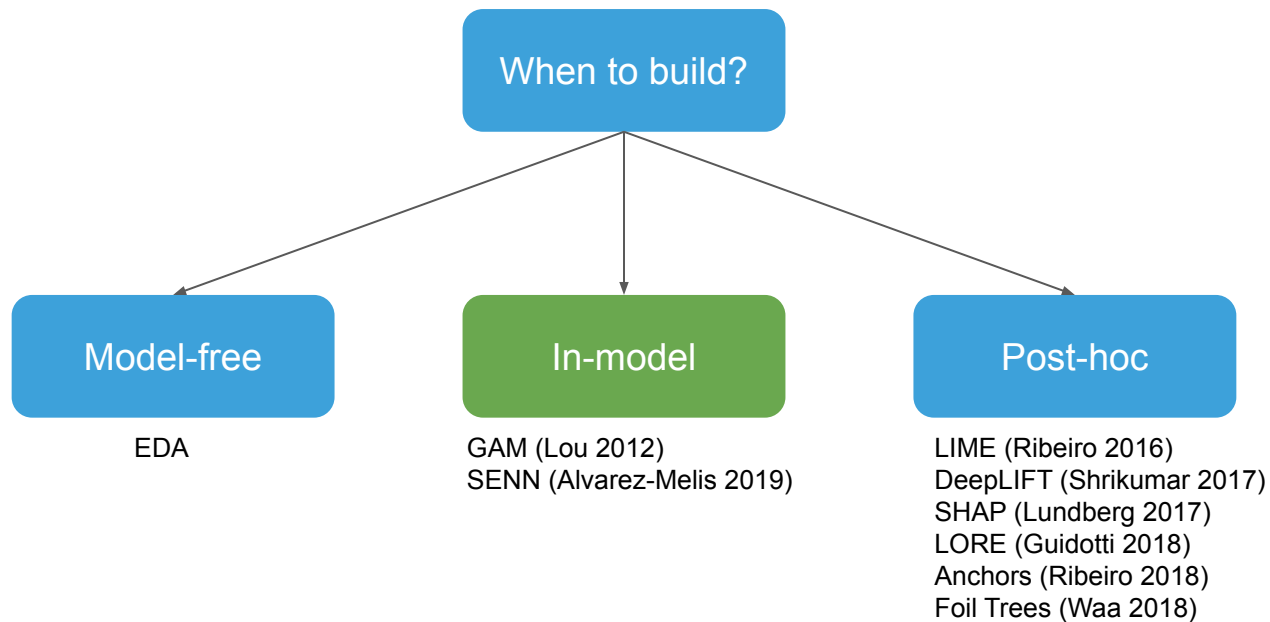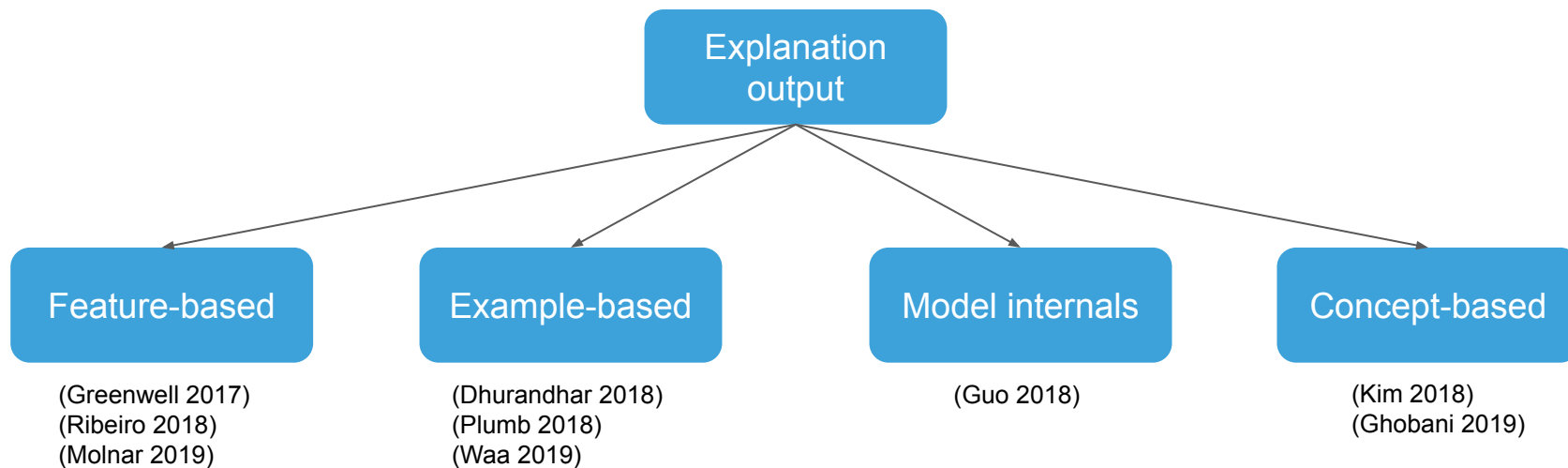- Efficiency
- Better & faster decisions

- Efficiency
- Iterate and debug models

- Improve outcomes
- Reduce friction while feeling safe

- Audit and Assess
- if the system is compliant.

*The transaction is **suspicious** because MCC = 7801.*

*The transaction is **suspicious** because it contains Suspicious Items.*

| Motivation | **Related Work** | Solution | Experiment | Conclusion |

# XAI Personas - Different personas, different XAI needs...

**Fields of expertise:**

- Domain knowledge
- No ML knowledge

- Data Science + ML

- Know what they want
- No domain nor ML knowledge

- Regulations/law
- Limited domain & ML knowledge

### Human-in-the-loop

### Data Scientist

### Decision Subject

### Regulator

**Goals:**

- Efficiency
- Better & faster decisions

- Efficiency
- Iterate and debug models
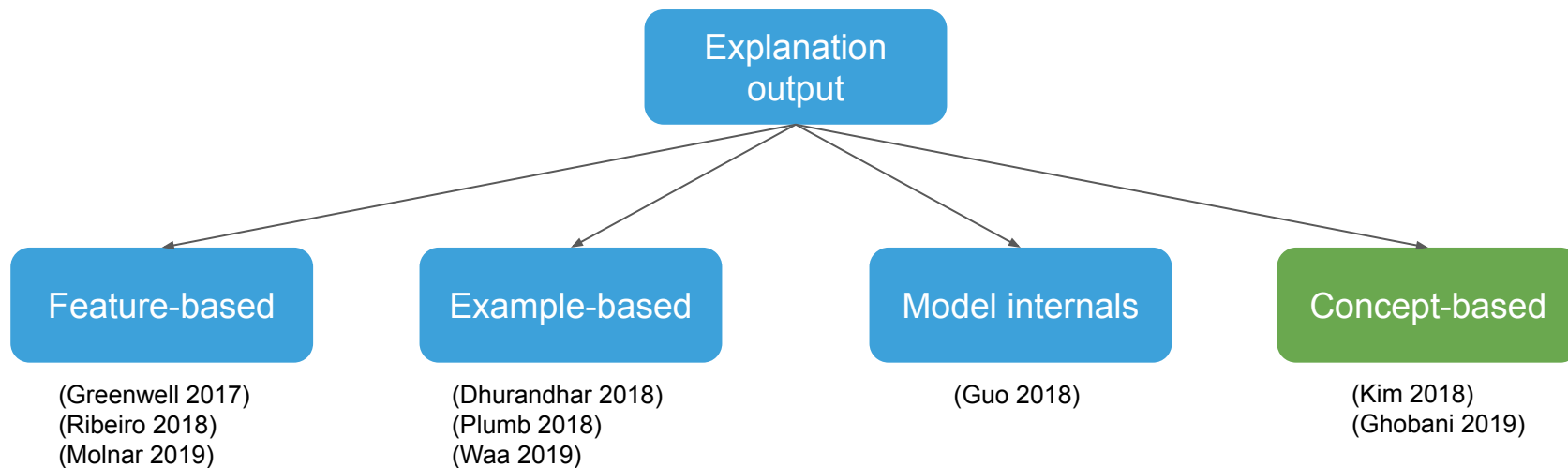
- Improve outcomes
- Reduce friction while feeling safe

- Audit and Assess
- if the system is compliant.

# The optimal choice of explanation depends on the end persona!

| Motivation | **Related Work** | Solution | Experiment | Conclusion |

# Interpretability: When to build

When to build?

Model-free

In-model

Post-hoc

EDA

GAM (Lou 2012)
SENN (Alvarez-Melis 2019)

LIME (Ribeiro 2016)
DeepLIFT (Shrikumar 2017)
SHAP (Lundberg 2017)
LORE (Guidotti 2018)
Anchors (Ribeiro 2018)
Foil Trees (Waa 2018)

Motivation | **Related Work** | Solution | Experiment | Conclusion

# Interpretability: When to build

When to build?

Model-free

In-model

Post-hoc

EDA

GAM (Lou 2012)
SENN (Alvarez-Melis 2019)

LIME (Ribeiro 2016)
DeepLIFT (Shrikumar 2017)
SHAP (Lundberg 2017)
LORE (Guidotti 2018)
Anchors (Ribeiro 2018)
Foil Trees (Waa 2018)

Motivation | **Related Work** | Solution | Experiment | Conclusion

# Taxonomy of Explanations' Output

Explanation output

Feature-based

Example-based

Model internals

Concept-based

(Greenwell 2017)
(Ribeiro 2018)
(Molnar 2019)

(Dhurandhar 2018)
(Plumb 2018)
(Waa 2019)

(Guo 2018)

(Kim 2018)
(Ghobani 2019)

| Motivation | Related Work | Solution | Experiment | Conclusion |

# Taxonomy of Explanations' Output

```
                          ┌─────────────────┐
                          │   Explanation   │
                          │     output      │
                          └─────────────────┘
```

| Feature-based | Example-based | Model internals | Concept-based |
|---|---|---|---|
| (Greenwell 2017) | (Dhurandhar 2018) | (Guo 2018) | (Kim 2018) |
| (Ribeiro 2018) | (Plumb 2018) | | (Ghobani 2019) |
| (Molnar 2019) | (Waa 2019) | | |

# There are many different methods and libraries available



feedzai

Alibi

TreeInterpreter

DeepLift

LORE

TCAV

SHAP

DiCE

SENN

LIME

Attention

ACE

Anchors

Saliency

GradCAM

GuidedBackprop

Integrated Gradients

And others...

- ACE is a global (and local), model-specific and concept-based explanation method that automatically groups input features into high-level concepts.

- The concepts are represented by groups of pixels (segments).



(a) Multi-resolution segmentation of images    (b) Clustering similar segments and removing outliers    (c) Computing saliency of concepts

Importance Scores

0.8

0.7

0.4

# Our Explainability Requirements

To best to our knowledge, there is state-of-the-art XAI method that satisfies our explainability requirements.

| In-Model | Local | Concept-based | Tabular data |

# Proposed Solution

# Jointly learned cOncept-based ExpLanations (JOEL)

- JOEL, a NN-based framework to jointly learn a decision-making task and associated domain knowledge explanations.

# Jointly learned cOncept-based ExpLanations (JOEL)

- JOEL, a NN-based framework to jointly learn a decision-making task and associated domain knowledge explanations.

- JOEL is a self-explainable model, i.e., it incorporates the interpretability architecturally, allowing to produce the decision and also the explanations related to its decision.

# Jointly learned cOncept-based ExpLanations (JOEL)

feedzai

- JOEL, a NN-based framework to jointly learn a decision-making task and associated domain knowledge explanations.

- JOEL is a self-explainable model, i.e., it incorporates the interpretability architecturally, allowing to produce the decision and also the explanations related to its decision.

- JOEL provides high-level insights about the model's predictions that very much resemble the domain experts' own reasoning.

# JOEL architecture

# JOEL architecture



© 2021 Feedzai. This presentation is proprietary.

# JOEL architecture



$$\mathcal{L}(\boldsymbol{x}, \boldsymbol{y}_E, \boldsymbol{y}_D) = \mathcal{L}_D(\hat{\boldsymbol{y}}_D, \boldsymbol{y}_D) + \mathcal{L}_E(\hat{\boldsymbol{y}}_E, \boldsymbol{y}_E)$$

| Motivation | Related Work | **Solution** | Experiment | Conclusion |

# Concept-Based Explainability and its limitations

Our problem is characterized by having:

- High-resources for the decision task **but** low-resources for the explainability task;

- Out-of-the-shelf domain knowledge (with no added cost).

Research Question:

*Can we **do better** than the **fully supervised** and **low-resources** baseline?*

# Challenges

## Label Scarcity
(insufficient concept labels)

## Multi-task Learning
(how to explain the model's predictions?)

- Good DL generalization requires massive datasets;

- Labeling campaigns are arduous and expensive to carry;

- Explainability task should explain the decision task;

- Explanations must reflect the human-in-the-loop's reasoning.

# Possible solutions

1. Explore **Weak Supervision** techniques;

    - Can we leverage domain expertise and already existing components in Human-AI system?

| Motivation | Related Work | **Solution** | Experiment | Conclusion |

# Possible solutions

1. Explore **Weak Supervision** techniques;

    - Can we leverage domain expertise and already existing components in Human-AI system?

2. Explore different **Learning Strategies**;

    - Use noisy labels only?
    - Use noisy labels and then fine-tune using golden labels?
    - Mix both labels?

# Implementation Workflow



Label Scarcity

Multi-task learning

# Implementation Workflow



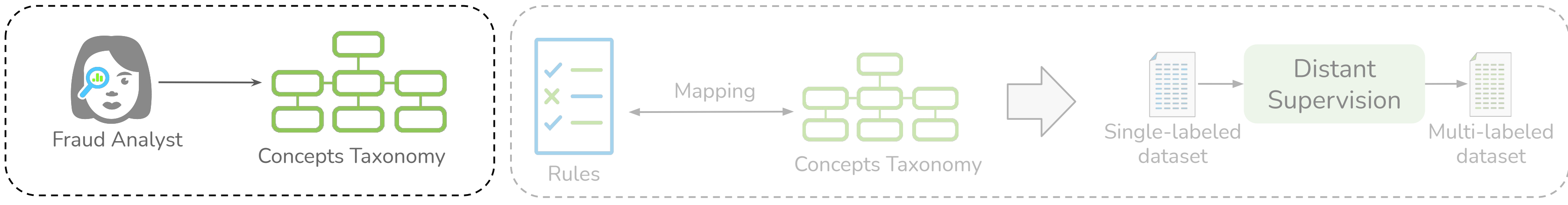


1 **Domain Expert defines concepts that will be used as explanations.**

2 **Apply Distant Supervision.**

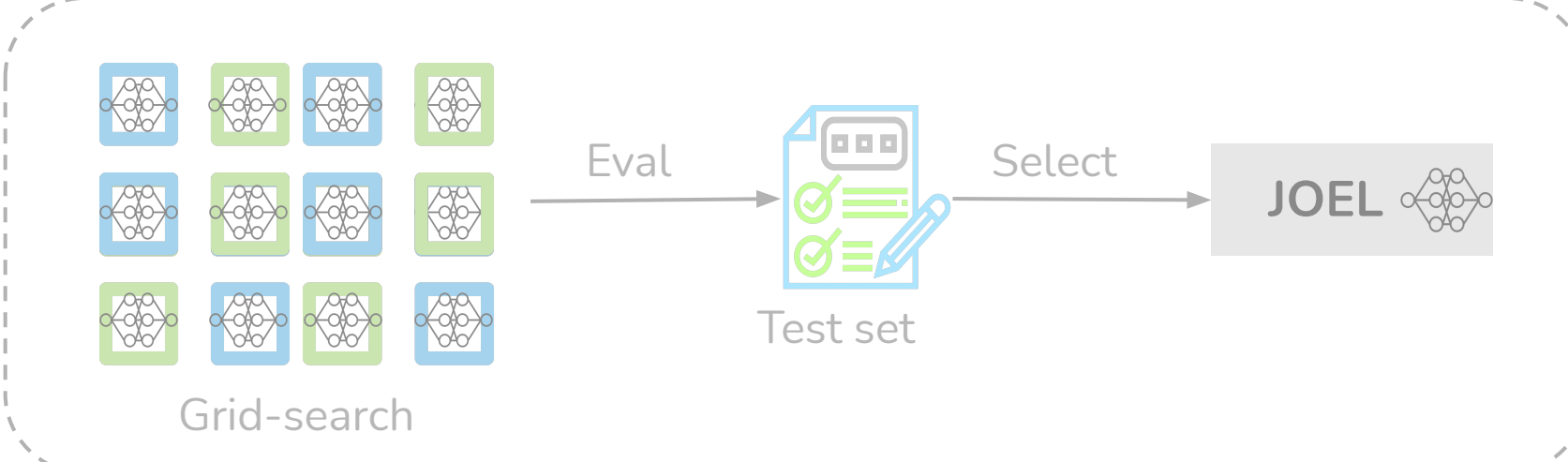


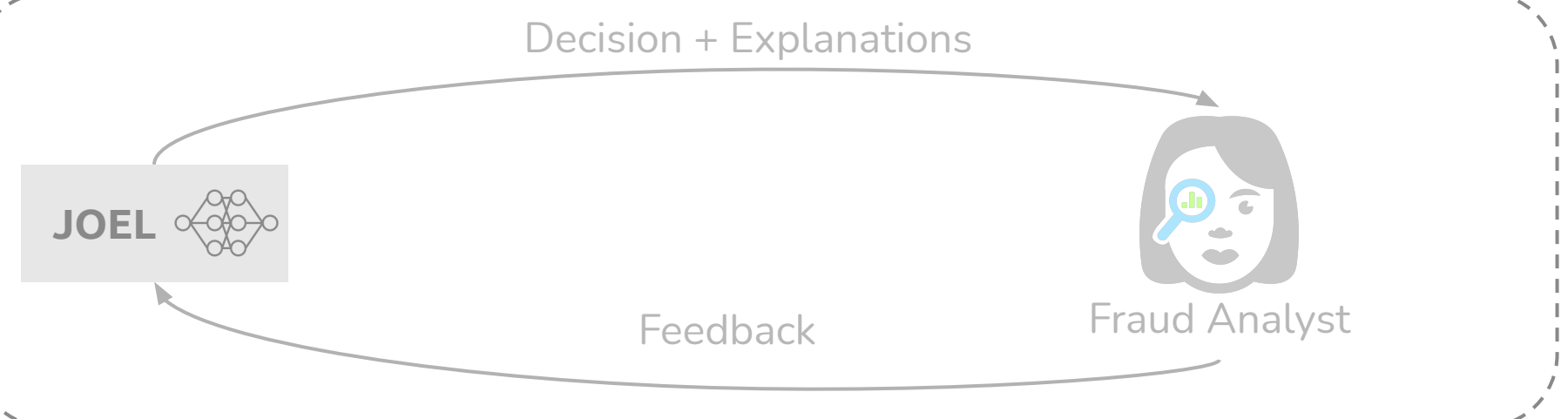


3 **Run a grid search to find the best hyperparameters for Distant Supervision and Supervised Learning.**

4 **JOEL's online loop generates concept-based explanations and collects human feedback in a fraud detection setting.**

# Fraud Taxonomy: Fraudulent Concepts

Suspicious device

High speed ordering

Suspicious items

Suspicious delivery

Suspicious billing shipping

Suspicious email

Suspicious customer

Suspicious IP

Other Fraud

Suspicious payment

Nothing suspicious

All or most details match

Good customer history

Other Legit

**9** Fraudulent concepts (+ **1** Other)

| Motivation | Related Work | Solution | **Experiment** | Conclusion |

# Fraud Taxonomy: Legitimate Concepts

feedzai

Suspicious device

High speed ordering

Suspicious items

Suspicious delivery

Suspicious billing shipping

Suspicious email

Suspicious customer

Suspicious IP

Other Fraud

Suspicious payment

Nothing suspicious

Good customer history

All or most details match

Other Legit

**3** Legitimate concepts (+ **1** Other)

Motivation | Related Work | Solution | **Experiment** | Conclusion

# Implementation Workflow



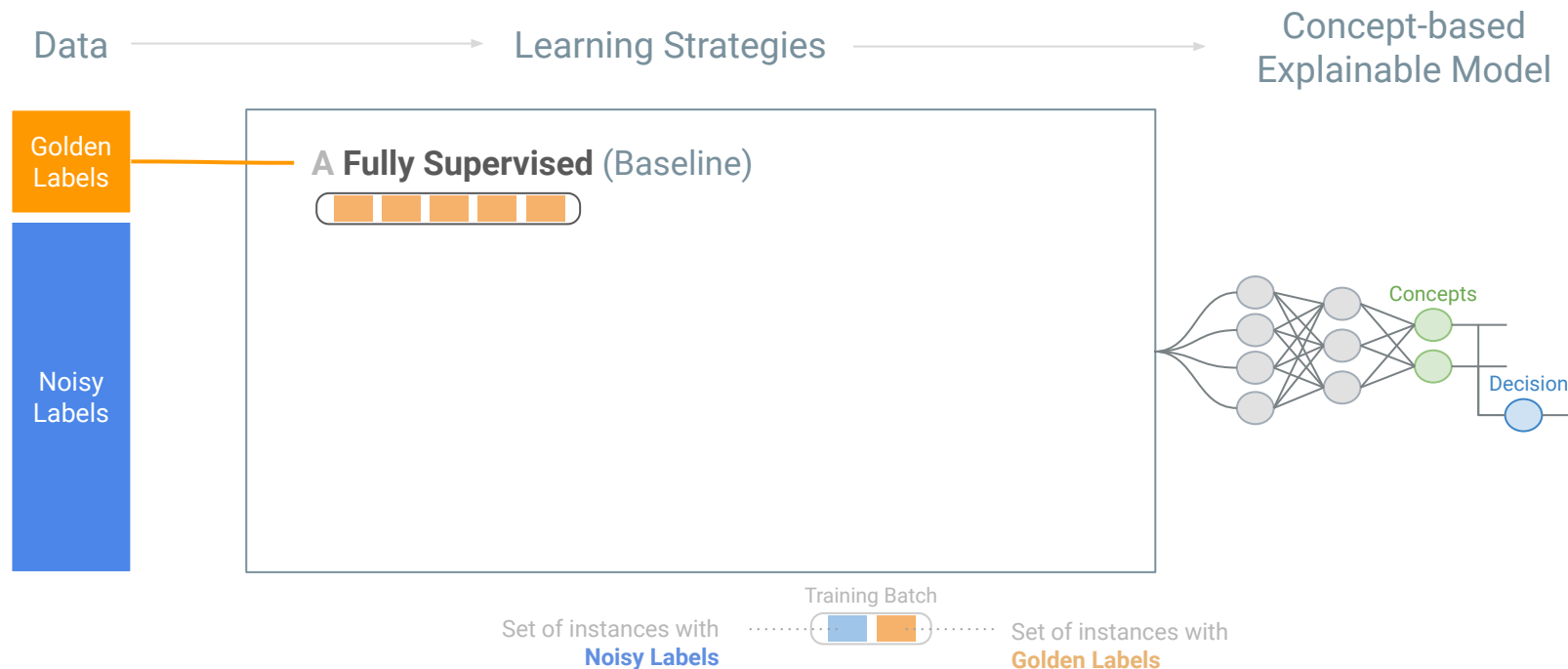1. **Domain Expert defines concepts that will be used as explanations.**
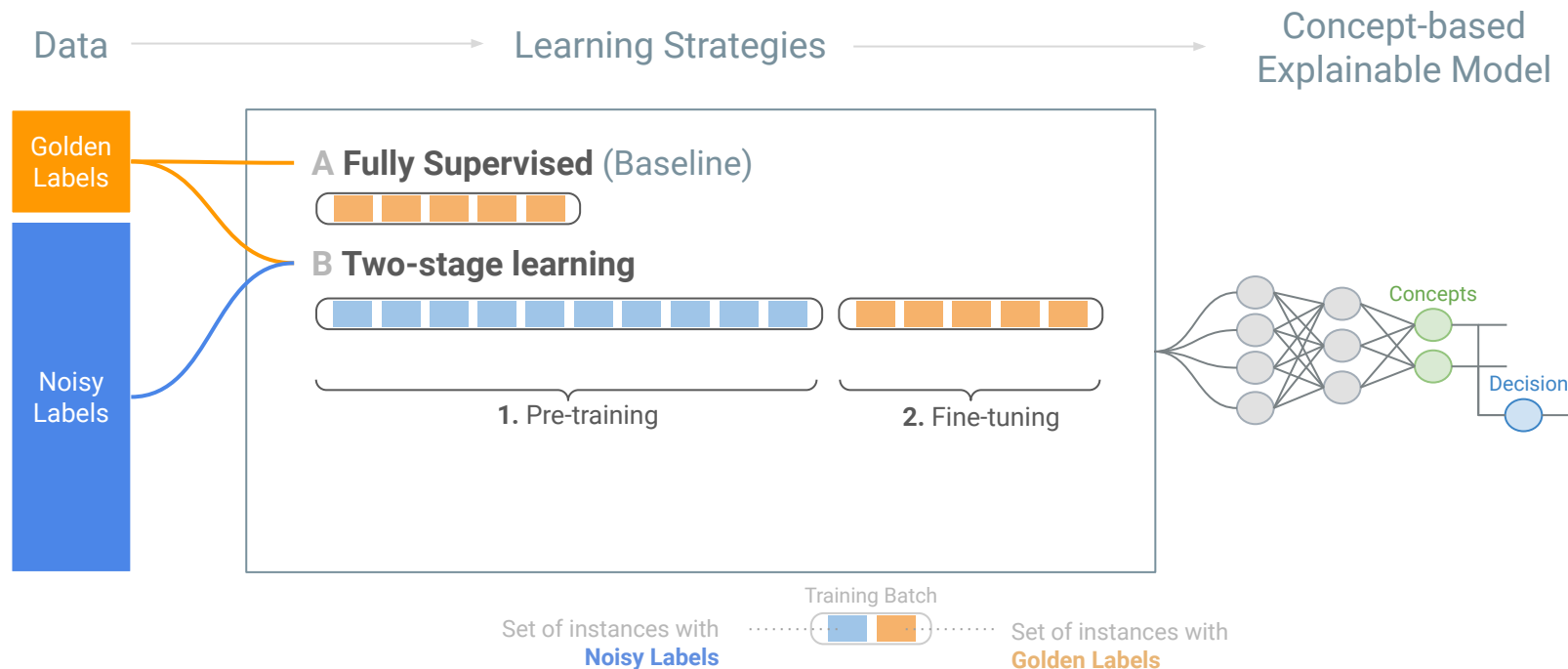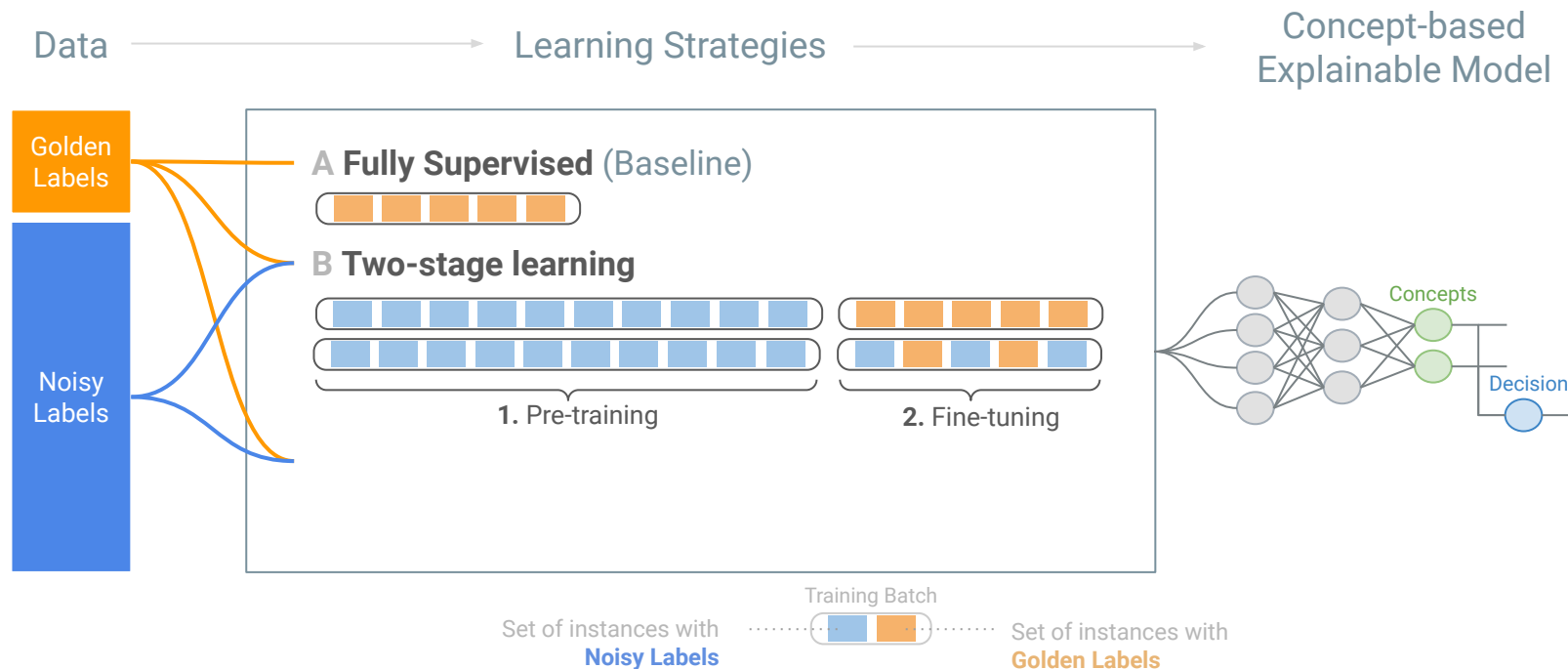
2. **Apply Distant Supervision.**

3. **Run a grid search to find the best hyperparameters for Distant Supervision and Supervised Learning.**

4. **JOEL's online loop generates concept-based explanations and collects human feedback in a fraud detection setting.**

| Motivation | Related Work | **Solution** | Experiment | Conclusion |

# Tackle Label Scarcity through Weak Supervision

1. Apply a **Distant Supervision** technique using available domain knowledge.



Rules

Unlabeled Data
(no concept labels)

*Rule-to-Concept* Mapping

| Rule description | Mapped concepts |
|---|---|
| Order contains risky product styles. | Suspicious Items |
| User tried *n* different cards last week. | Suspicious Customer, Suspicious Payment |

Noisy Labeled Data
(weak concept labels)

Concepts Taxonomy

| Motivation | Related Work | **Solution** | Experiment | Conclusion |
|---|---|---|---|---|

# Rules-Concepts Mapping for Distant Supervision

- One rule can be associated with one or more fraud concepts.

# Rules-Concepts Mapping for Distant Supervision

- One rule can be associated with one or more fraud concepts.

- By having a Rules-Concepts mapping, we can use "noisy labels" (mapped concepts) for each transaction using the rules that were triggered.

# Rules-Concepts Mapping for Distant Supervision

feedzai

- One rule can be associated with one or more fraud concepts.

- By having a Rules-Concepts mapping, we can use "noisy labels" (mapped concepts) for each transaction using the rules that were triggered.

| Canonical rule name | Description | Associated concepts |
|---|---|---|
| *amount_rule* | Transaction amount larger than $N | Suspicious Items |
| *expl_shippingip_mismatch* | Shipping and IP mismatch | Suspicious IP, Suspicious billing shipping |
| *jp__card_id_count_N__K_m* | The payment card id was detected more than N times in less than K minutes. | Suspicious Payment, High speed ordering |

# Implementation Workflow

① Domain Expert defines concepts that will be used as explanations.

② Apply Distant Supervision.

③ Run a grid search to find the best hyperparameters for Distant Supervision and Supervised Learning.

④ JOEL's online loop generates concept-based explanations and collects human feedback in a fraud detection setting.
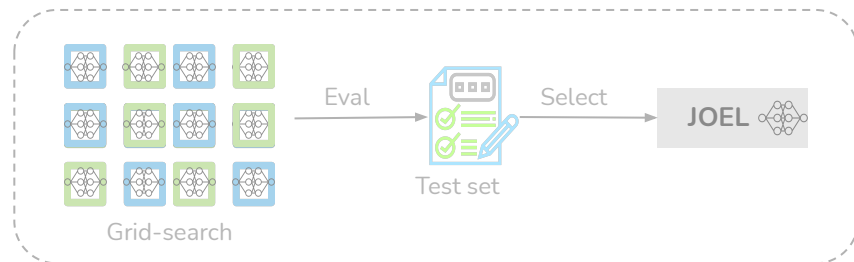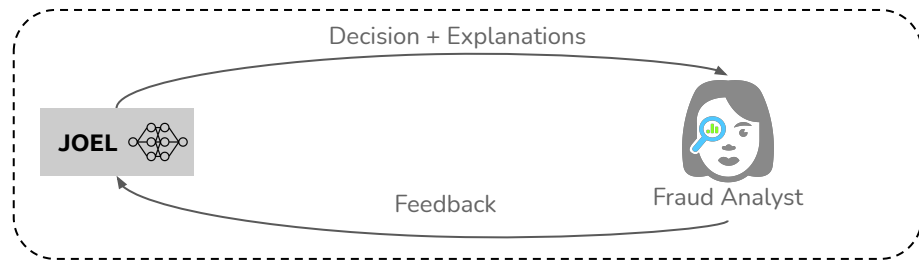
# Jointly learn to explain and decide: Learning Strategies



Data → Learning Strategies → Concept-based Explainable Model

Golden Labels

Noisy Labels

**A Fully Supervised** (Baseline)

Concepts

Decision

Training Batch

Set of instances with **Noisy Labels**

Set of instances with **Golden Labels**

# Jointly learn to explain and decide: Learning Strategies

# Jointly learn to explain and decide: Learning Strategies



Data → Learning Strategies → Concept-based Explainable Model

Golden Labels

Noisy Labels

**A** **Fully Supervised** (Baseline)

**B** **Two-stage learning**

**1.** Pre-training

**2.** Fine-tuning

Concepts

Decision

Training Batch

Set of instances with **Noisy Labels**

Set of instances with **Golden Labels**

# Jointly learn to explain and decide: Learning Strategies



Data → Learning Strategies → Concept-based Explainable Model

Golden Labels

Noisy Labels

**A** **Fully Supervised** (Baseline)

**B** **Two-stage learning**

**1.** Pre-training    **2.** Fine-tuning

**C** **Hybrid learning**

Concepts

Decision

Training Batch

Set of instances with **Noisy Labels**    Set of instances with **Golden Labels**

| Motivation | Related Work | **Solution** | Experiment | Conclusion |

# Implementation Workflow



1. Domain Expert defines concepts that will be used as explanations.

2. Apply Distant Supervision.

3. Run a grid search to find the best hyperparameters for Distant Supervision and Supervised Learning.

4. JOEL's online loop generates concept-based explanations and collects human feedback in a fraud detection setting.

# Human Feedback

- JOEL's explanations and prediction are shown to analyst through Web UI;

- A human-in-the-loop (*e.g.,* fraud analyst) makes a decision based on this information;

- When submitting its review, it also gives feedback about the concepts that led to his decision (and also about the decision task);

- This feedback can be used to continuously improve predictive accuracy, and also explainability of the model.

# Experiments

# Experimental Setup

Payment retailer aims to deploy multi-task approach for concept-based explanations.

- **Binary decision task**: detect fraudulent transactions;

- **Multi-Label explainability task** (with 14 concepts): output high-level explanation about the model's prediction.

| Labels | Availability |
| --- | --- |
| Golden decision | High (~6M, ~2% fraud rate) |
| Noisy Explainability | High (~6M, ~2% fraud rate) |
| Golden* explainability | Low (~1.3k, ~37% fraud rate) |

# Evaluation Metrics

## Decision Task

Fraud recall @ 5% FPR



ROCs on Test

- Standard FF
- JOEL (best on decision)
- LightGBM
- JOEL (best on explainability)

## Explainability Task

Mean Average Precision (mAP)

$$AP = \frac{\sum_{k=1}^{n}(P(k) * rel(k))}{number\ of\ relevant\ items}$$

$$MAP = \frac{1}{Q}\sum_{q=1}^{Q} AP(q)$$

*Can we **do better** than the **fully supervised** and **low-resources** baseline?*

Model development:

- Run the same hyperparameter grid for each variant;
  (*i.e.*, number and dimension of hidden layers, learning rate, explainability task importance)


- Run 2 random seeds.

# Fully Supervised (baseline) - Test set

feedzai

**Explainability Task (Mean Average Precision)**



**Decision Task (Fraud Recall)**

Full supervised learning on a small dataset yields **poor results** in the decision task.

| Motivation | Related Work | Solution | **Experiment** | Conclusion |

# Fully Supervised (baseline) - Test set

feedzai

**Explainability Task (Mean Average Precision)**



**Decision Task (Fraud Recall)**

| | Full Supervision |
|---|---|
| Recall | [0.04; 0.15] |
| mAP | [0.52; 0.54] |

| Motivation | Related Work | Solution | **Experiment** | Conclusion |

# Two-stage learning - Test set

feedzai

**First stage:**
Train base models using
Distant Supervision (Noisy Labels).

**Second stage:**
Fine-tune base models with
Golden Labels.

**Main goal**:
improve **explainability** task
without hurting the **decision** task.

| | Full Supervision | Two-stage (base models) | Two-stage (w/ fine tuning) |
|---|---|---|---|
| Recall | [0.04; 0.15] | | |
| mAP | [0.52; 0.54] | | |

# Two-stage learning - Test set

Base models (random seeds 10 and 42)



**Explainability Task (Mean Average Precision)** vs **Decision Task (Fraud Recall)**

|  | **Full Supervision** |
|---|---|
| Recall | [0.04; 0.15] |
| mAP | [0.52; 0.54] |

- Explainability performance <span style="color:darkred">degrades</span> when using the larger but noisy explainability dataset.

- Training with larger dataset <span style="color:green">improves</span> decision task;

# Two-stage learning - Test set

Base models (random seeds 10 and 42)

**Explainability Task (Mean Average Precision)**



Decision Task (Fraud Recall)

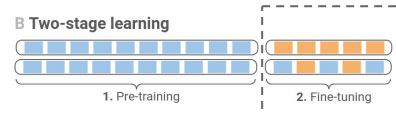|  | **Full Supervision** | **Two-stage (base models)** |
|---|---|---|
| Recall | [0.04; 0.15] | [0.6; 0.78] |
| mAP | [0.52; 0.54] | [0.24; 0.29] |

# Two-stage learning - Test set

**Second stage:**
Fine-tune base models with Golden Labels.

**Main goal**:
improve **explainability** task without hurting the **decision** task.

|  | Full Supervision | Two-stage (base models) | Two-stage (w/ fine tuning) |
|---|---|---|---|
| Recall | [0.04; 0.15] | [0.6; 0.78] | |
| mAP | [0.52; 0.54] | [0.24; 0.29] | |

# Two-stage learning - Test set

Pick **Pareto Optimal** models from Stage 1.

Fine-tune them with:

- Hyperparameters;
  (*e.g.*, learning rate, epochs, batch size)

- Loss scalers;
  (*e.g.*, 0.75, 0.5, 0.25)

- Freeze and unfreeze of layers;

- Different batch techniques.
  (*e.g.*, hybrid batching - add % of noisy labels in each batch)

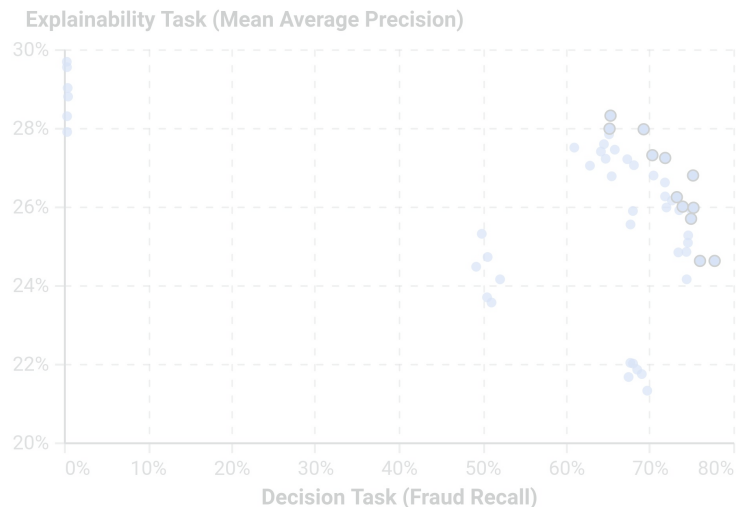**Second stage:**
Fine-tune base models with Golden Labels.

**Main goal**:
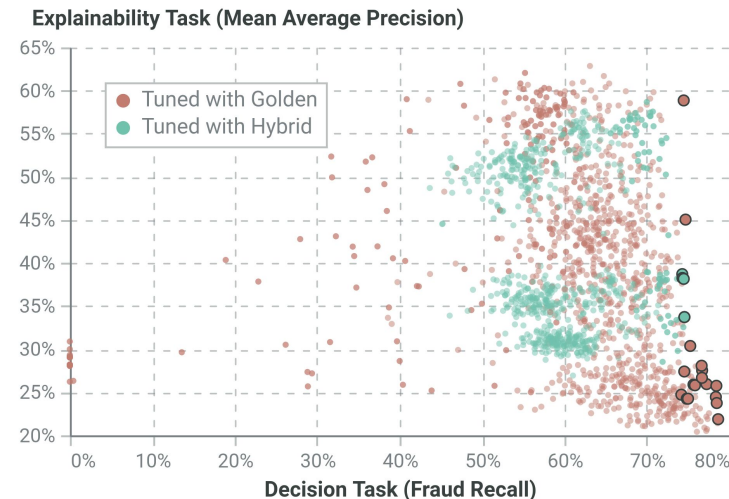improve **explainability** task without hurting the **decision** task.

# Two-stage learning - Test set
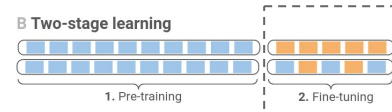
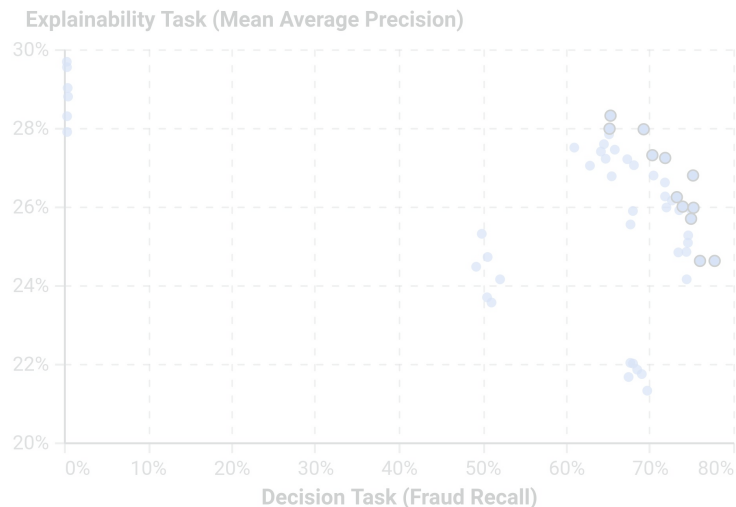Base models (random seeds 10 and 42)

Fine-tuned models



**Explainability Task (Mean Average Precision)**

Decision Task (Fraud Recall)



**Explainability Task (Mean Average Precision)**

- Tuned with Golden
- Tuned with Hybrid

Decision Task (Fraud Recall)

For same decision task performance, we can improve explainability task.

| | | | | |
|---|---|---|---|---|
| Motivation | Related Work | Solution | **Experiment** | Conclusion |

# Two-stage learning - Test set

## Base models (random seeds 10 and 42)

**Explainability Task (Mean Average Precision)**



Decision Task (Fraud Recall)

## Fine-tuned models

**Explainability Task (Mean Average Precision)**



- Tuned with Golden
- Tuned with Hybrid

Decision Task (Fraud Recall)

| | Full Supervision | Two-stage (base models) | Two-stage (w/ fine tuning) |
|---|---|---|---|
| Recall | [0.04; 0.15] | [0.6; 0.78] | [0.74; 0.78] |
| mAP | [0.52; 0.54] | [0.24; 0.29] | [0.24; **0.63**] |

feedzai

Each training batch contains a fraction of Explainability Golden Labels (Our experiments use 10%).

# Hybrid learning - Test set

**Explainability Task (Mean Average Precision)**



|  | Full Supervision | Two-stage (base models) | Two-stage (w/ fine tuning) | Hybrid |
|---|---|---|---|---|
| Recall | [0.04; 0.15] | [0.6; 0.78] | [0.74; 0.78] | **[0.65; 0.71]** |
| mAP | [0.52; 0.54] | [0.24; 0.29] | [0.24; 0.63] | **[0.4; 0.49]** |

Better than **fully supervised** in decision and **two-stage base** models in explainability, but seems to be worse than **two-stage fine-tuned**.

# All Pareto from learning strategies - test set

feedzai

**Explainability Task (Mean Average Precision)**



**Decision Task (Fraud Recall)**

- ● Fully Supervised (baseline)
- ● Two-stage (Tuned with Golden)
- ● Hybrid
- ● Two-stage (base models)
- ● Two-stage (Tuned with Hybrid)

**Preliminary results** seem to be **promising** but further experiments
(more seeds and more runs) to gain statistical confidence.

# Conclusions

# Recap

- Concept-based explainability through multi-task learning poses challenges:

  - Label scarcity;

  - Joint learning of decision and associated explanations.

- This work proposes to:

  - Use Distant Supervision and exploit the available off-the-shelf domain knowledge;

  - Use different Learning Strategies and combine label qualities to improve performance at both tasks.

# Conclusions

- The explanations should be tailored to the persona's knowledge and task performed;

  - Concept-based explanations are suitable to domain experts that make ML-informed decisions but lack ML knowledge!

- Experiment in a real-world e-commerce fraud detection dataset show:

  1. JOEL is able to learn both domain concept explanations and fraud decisions;

  2. Distant supervision allows us to overcome the label scarcity problem;

  3. There is no clear winner learning strategy.
     (it might depend on business requirements)

feedzai

# Questions?

vladimir.balayan@feedzai.com

catarina.belem@feedzai.com

a4338@fct.unl.pt

pedro.saleiro@feedzai.com

pedro.bizarro@feedzai.com