



# Trusted AI

## IBM Watson OpenScale

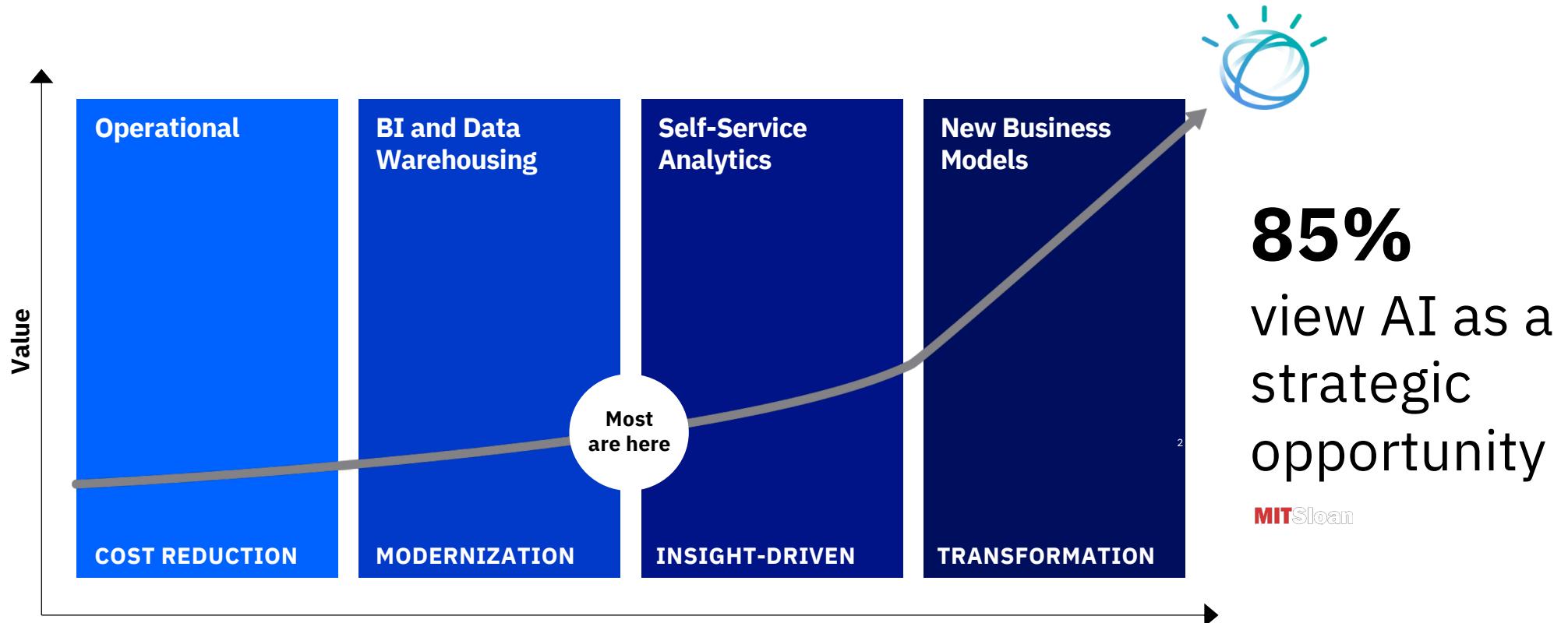
**Dr. Rudolf Pailer**  
Artificial Intelligence Practice Austria  
IBM Services

3. Oktober 2019

[linkedin.com/in/rudolfpailer](https://linkedin.com/in/rudolfpailer)



# I want AI !



# BUT..., business stakeholders do not trust AI.

60%

of companies see **regulatory constraints** as a barrier to implementing AI.

- IBM IBV AI 2018

63%

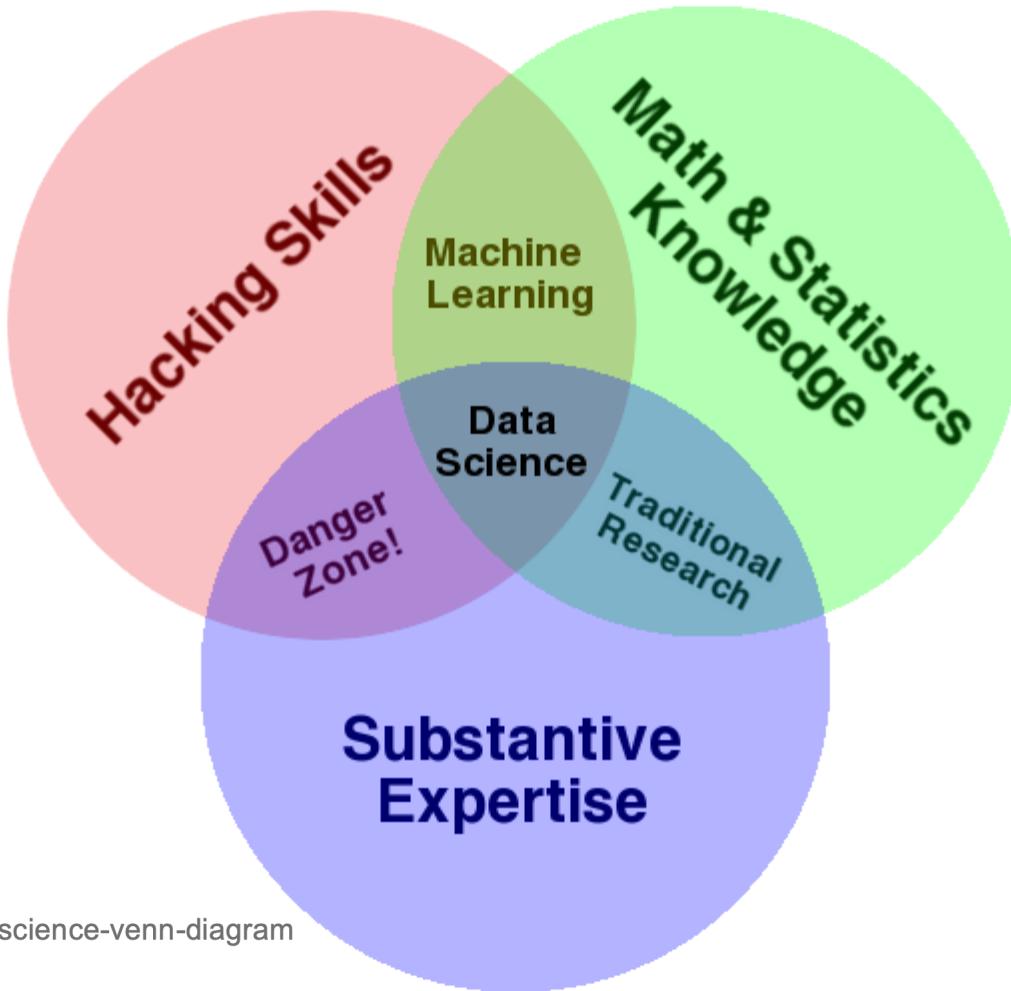
cite availability of **technical skills** as a challenge to implementation.

- IBM IBV AI 2018

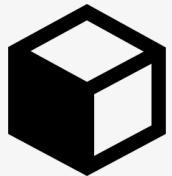
*Without expensive Data Science resources handholding multiple AI models in a production application:*

1. No way to **validate** if AI models are **compliant with regulations** and will achieve expected business outcomes before deploying
2. Difficult to **track and measure** indicators of business success in production
3. Resource intensive and unreliable processes for **ongoing business monitoring and compliance**
4. Impossible for business users to **feedback** subtle domain knowledge into model lifecycle

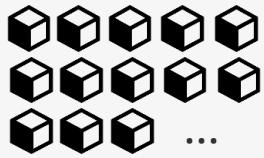
# Skill Requirements in Data Science & AI Projects



# Model management



Modeling & Evaluation



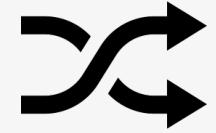
Model Versioning



Model Deployment



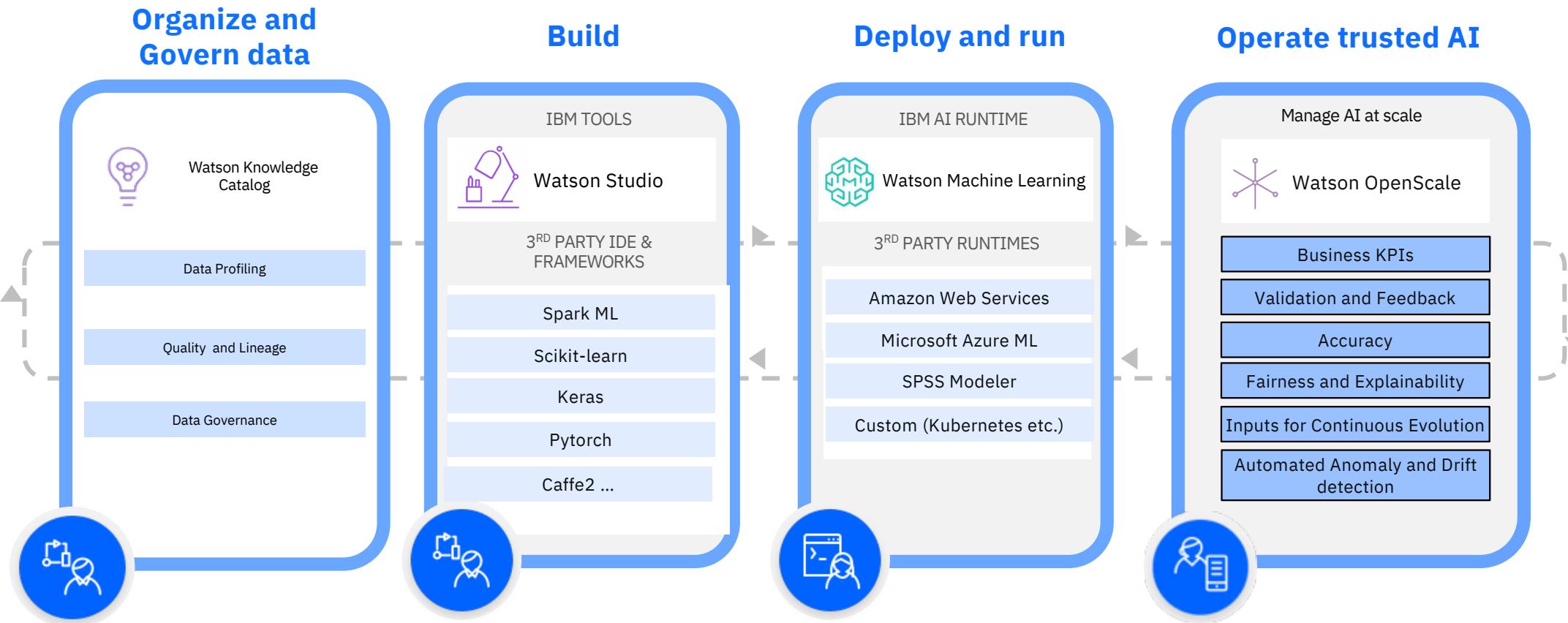
Model Monitoring



Dynamic Model Selection  
& Retraining

Data Science Solutions are **not** static by definition!

# Watson OpenScale along with Watson Studio and WML enables enterprises to operationalize AI across the enterprise



# Watson OpenScale will help validate and monitor AI models, deployed anywhere, to help comply with regulations and mitigate business risk

## Production monitoring for compliance and safeguards

Detect and mitigate model biases

Audit and Explain model decisions

Model Validation and acceptance

*Required in regulated industries and use cases – FSS, HR etc. in short term; others longer term\**

## Ensure that models are resilient to changing situations

Detect drift in data and anomaly in model behavior

Specific inputs and triggers to model lifecycle

*Required to meet transformational goals*

## Align model performance with business outcomes

Correlate model metrics and business KPIs to measure business impact

Actionable metrics and alerts

*Foundational to all AI implementations*

Current capability

Upcoming capability

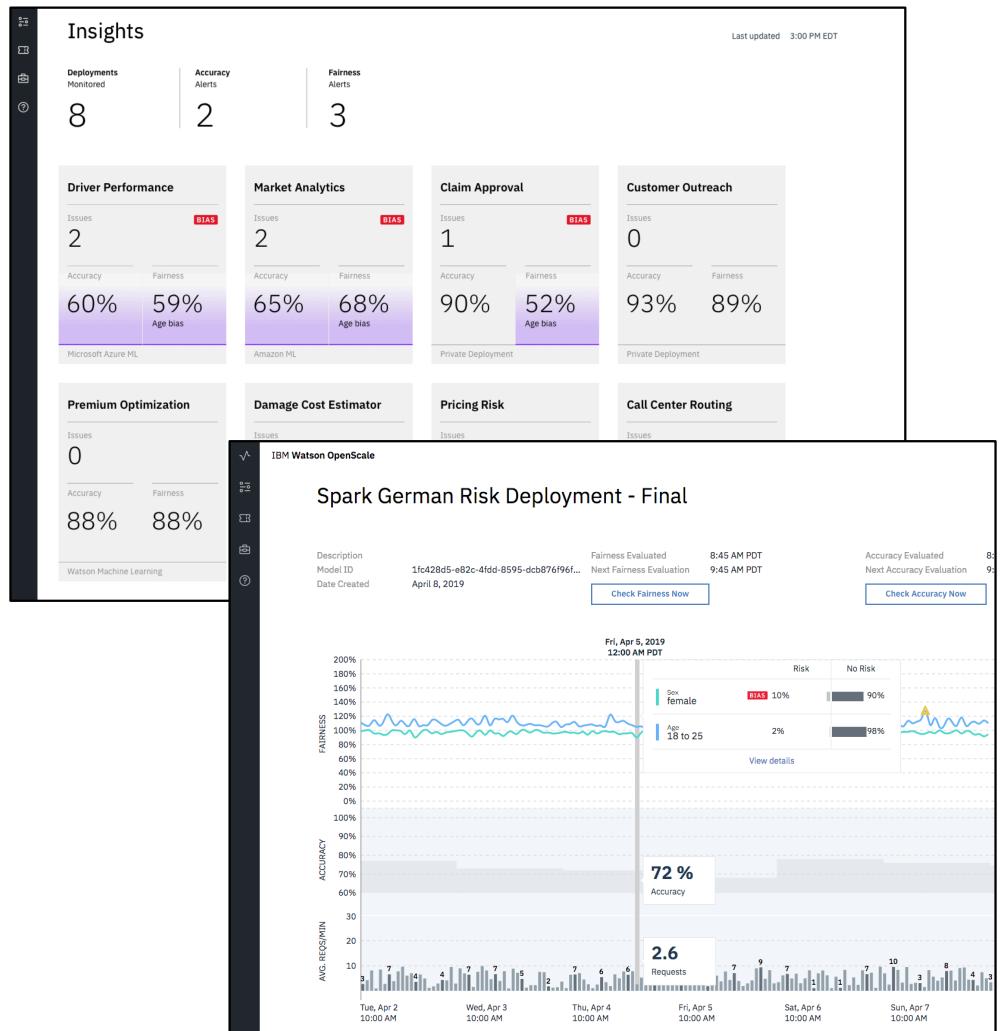
# Operations Dashboard

## Description:

Monitor deployed models in a single dashboard that can be filtered by deployment making it easy to manage AI in apps

## Value:

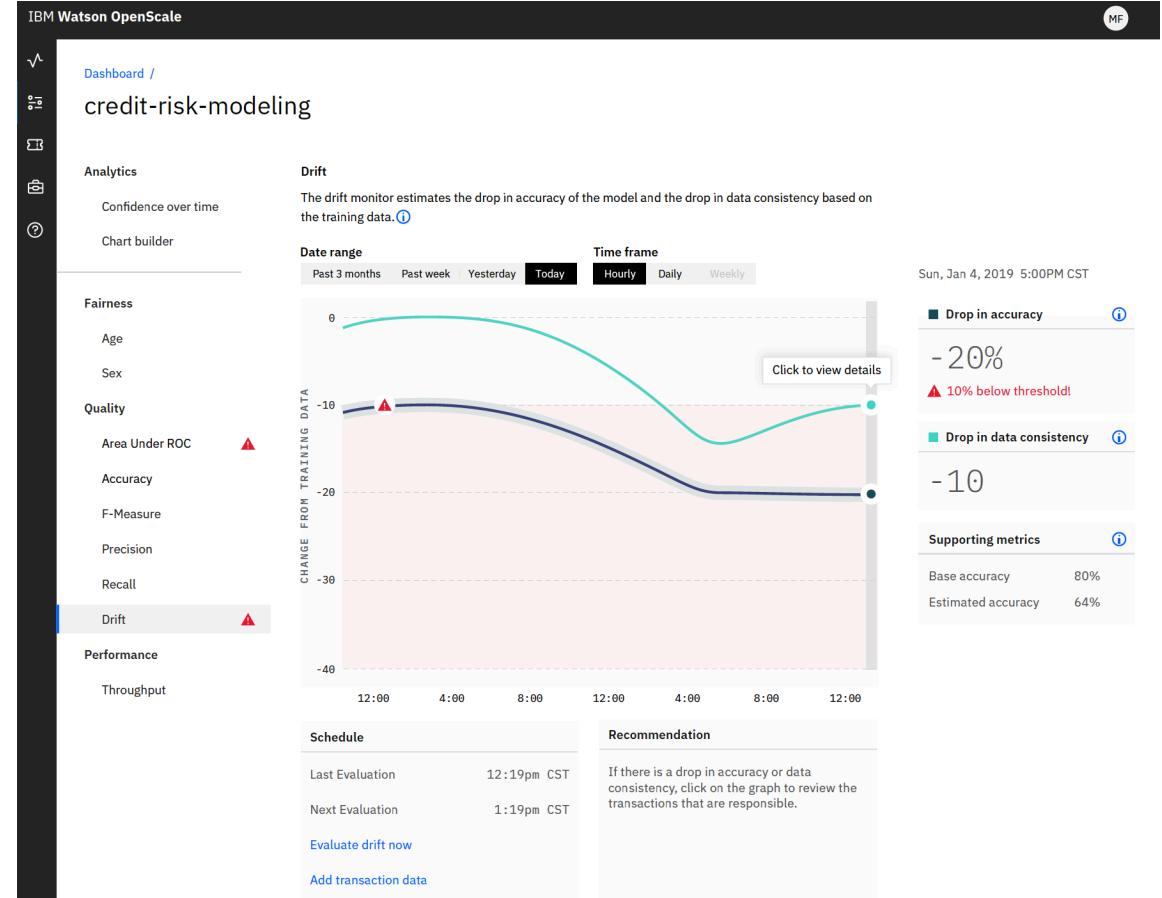
- Configure alerts or actions to be triggered when KPIs exceed threshold, ensuring model quality for improve business outcomes
- Measure model accuracy as it pertains to it's ability to deliver outcomes more accurate than knowledge workers



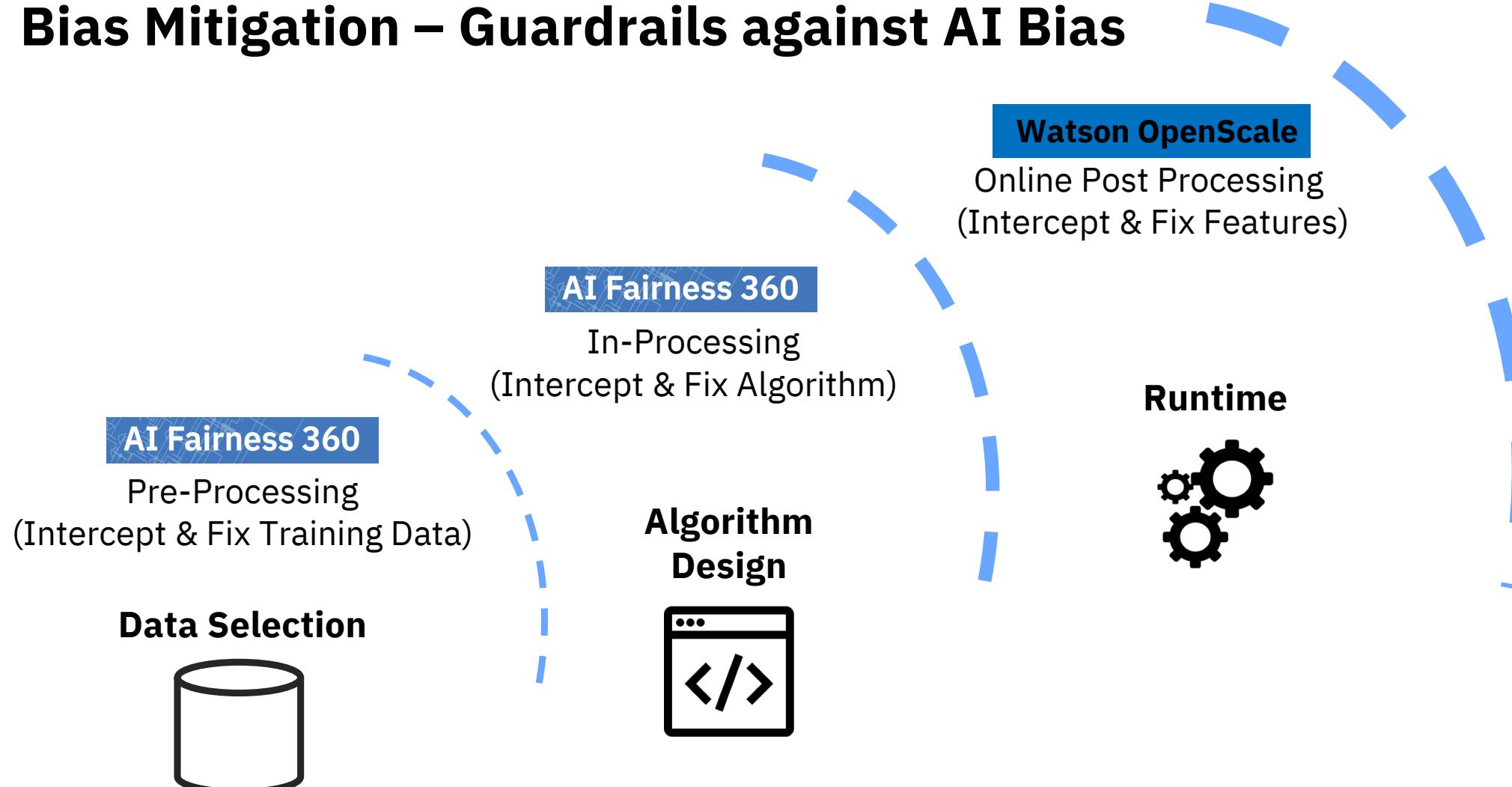
# Drift Detection in OpenScale

Drift Monitor in OpenScale measures two types of drifts:

- **Drop in accuracy:** It estimates the drop in accuracy of the model at runtime. The model accuracy could drop if there is an increase in transactions similar to those which the model was unable to evaluate correctly in the training data.
- **Drop in data consistency:** It estimates the drop in consistency of the data at runtime as compared to the characteristics of the data at training time.



# Bias Mitigation – Guardrails against AI Bias



# Ethics & Bias, COMPAS Recidivism

- Machine Bias: Pro Publica

“There’s software used across the country to predict future criminals. And it’s biased against blacks.”

<https://www.propublica.org/article/machine-bias-risk-assessments-in-criminal-sentencing>

- <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC5777393/>

“...We further show that a simple linear predictor provided with only two features is nearly equivalent to COMPAS with its 137 features.”



[Tweet](#)



Andrew Ng   
@AndrewYNg

AI+ethics is important, but has been partly hijacked by the AGI (artificial general intelligence) hype. Let's cut out the AGI nonsense and spend more time on the urgent problems: Job loss/stagnant wages, undermining democracy, discrimination/bias, wealth inequality.

[Original \(Englisch\) übersetzen](#)

6:01 nachm. · 11 Juni 18 aus Stanford, CA

2 469 Retweets 6 426 „Gefällt mir“-Angaben



Matt Leach @nextcontext · 11. Juni

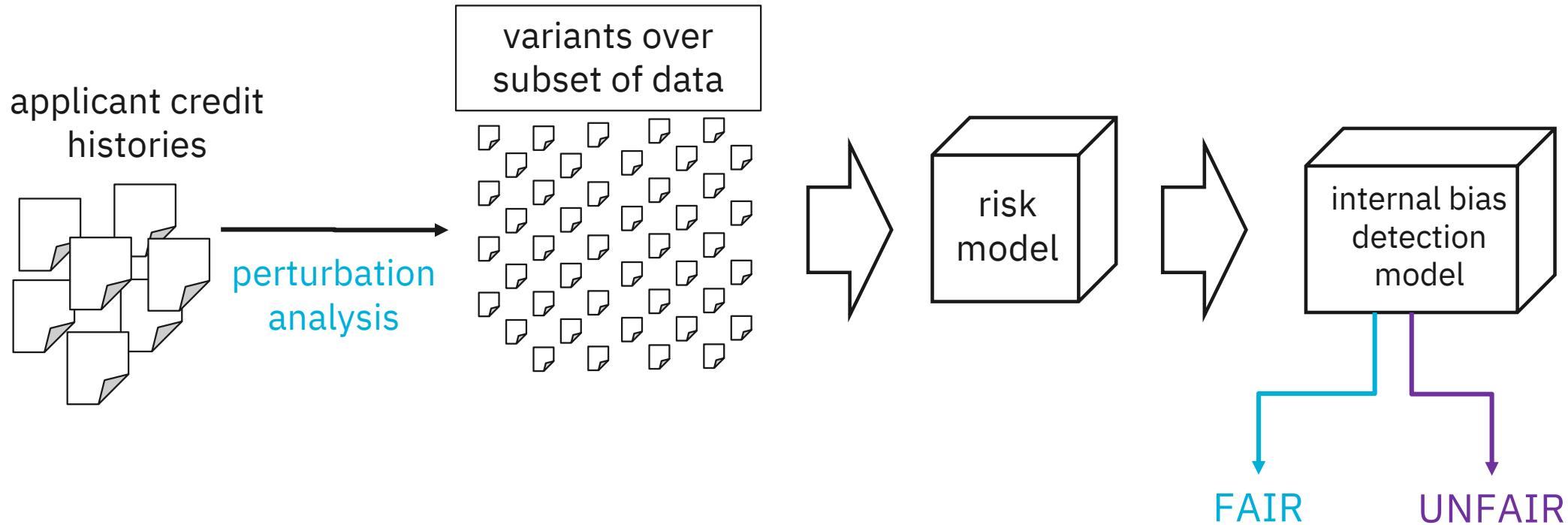
Antwort an @AndrewYNg

I think we need to be talking about both sets of risks, right now. The social/

Deine Antwort twittern

# How does Watson OpenScale mitigate fairness issues?

Calculated hourly over a sliding window



BIAS MITIGATION POST-PROCESSING FOR INDIVIDUAL AND GROUP FAIRNESS, Lohia et al.

<https://arxiv.org/abs/1812.06135>



# Model Fairness

## Description:

Models in production need to make fair decisions and can not be biased in their recommendations

## How Does it Work?

- Outcomes are selected as “favorable or unfavorable”
- “Favored Populations” and “protected populations” are selected where majority and minority groups are found
- A score is calculated based on the probability of favorable outcome for minority vs. probability of favorable outcome for majority

## Value Add

- Configure fairness metrics for continuous evolution – this provides quality checks so no biased data or models can drive unfair decisions in business applications
- Visualize the fairness of deployed models on the OpenScale dashboard, so model builders can act swiftly to mitigate bias in production deployments

The image displays two screenshots of the IBM Watson OpenScale interface. The top screenshot shows the 'Select the features to monitor' step. It lists various model features in a grid: CheckingStatus, LoanDuration, CreditHistory, LoanPurpose, LoanAmount, ExistingSavings, EmploymentDuration, InstallmentPercent, Sex, OthersOnLoan, CurrentResidenceDuration, OwnsProperty, Age, InstallmentPlans, and Housing. The 'Sex' feature is highlighted with a blue border. The bottom screenshot shows the 'Specify the favorable outcomes' step. It asks for 'Favorable values' and 'Unfavorable values'. Under 'Favorable values', there is a list with 'No Risk' and an 'Add' button. Under 'Unfavorable values', there is a list with 'Risk' and an 'Add' button.

# Bias Mitigation – Original Output

IBM AI

MF

## Claim Approval : Output

Original Output

De-biased Output

Age

October 14, 2018

08:00

AM

Only

52% of the group 18 to 23  
received the outcome status=Approved

By comparison

91% of the group 31 to 35  
received the outcome status=Approved

Accuracy

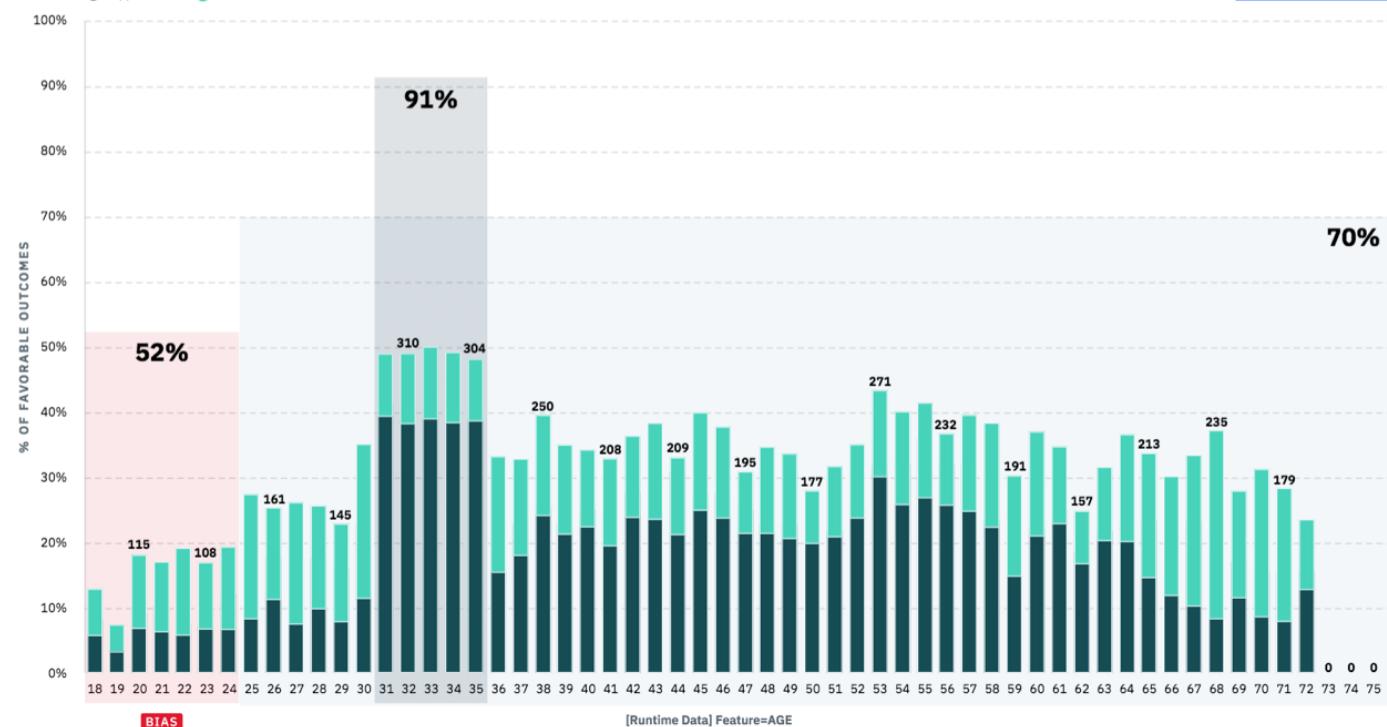
88%  
before de-bias

Model Output

● Approved    ● Denied

Runtime Data  Training Data

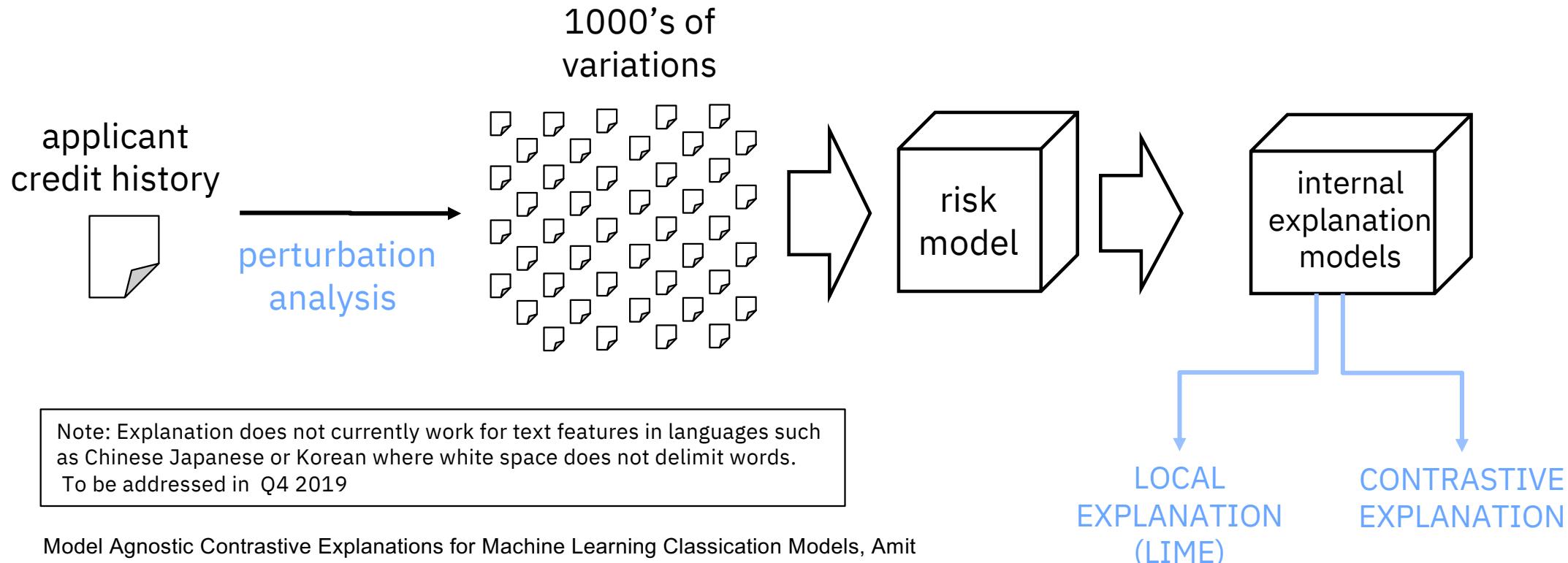
[View Transactions](#)



IBM

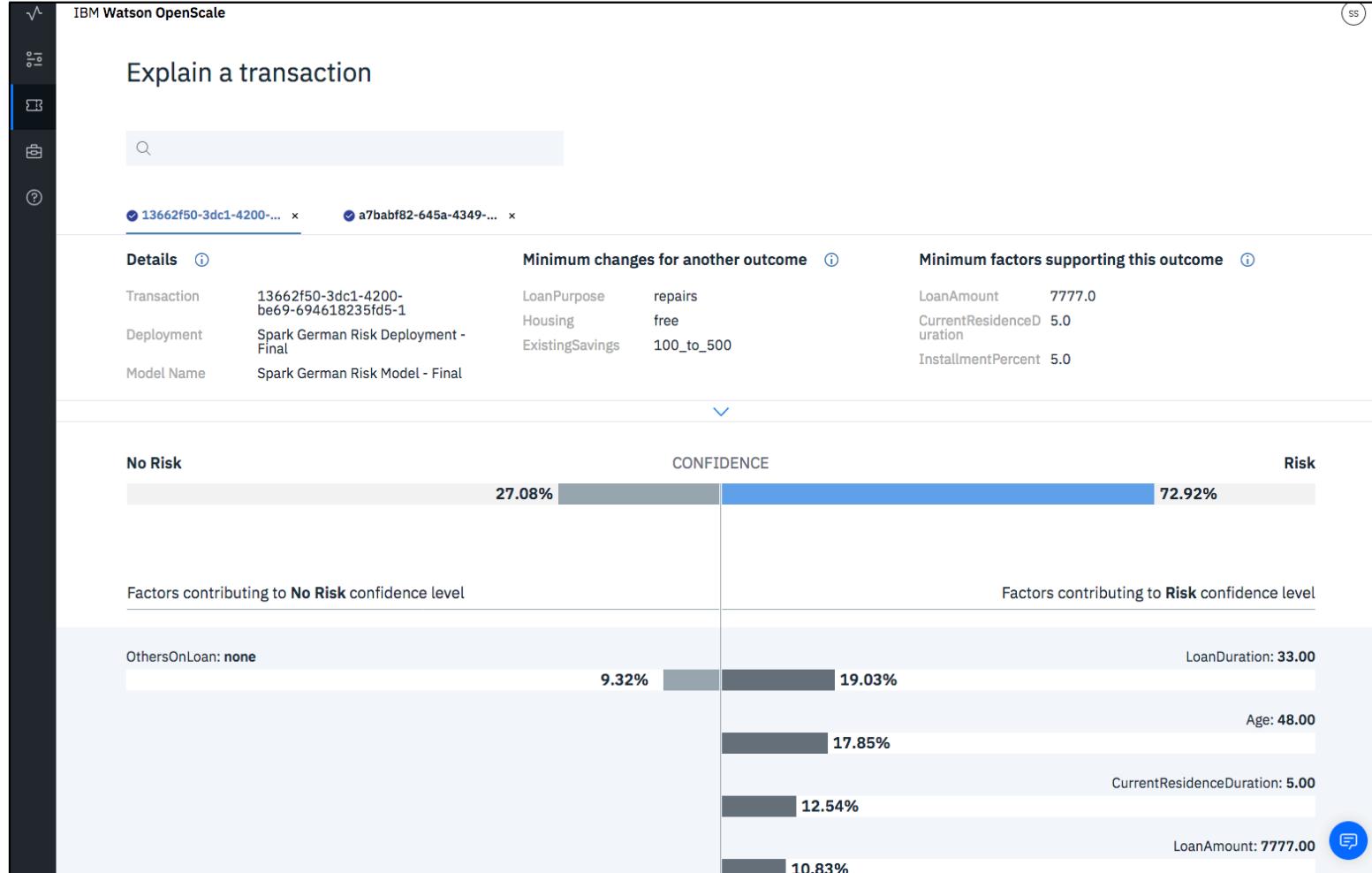
# How does Watson OpenScale explain a prediction?

Calculated upon request



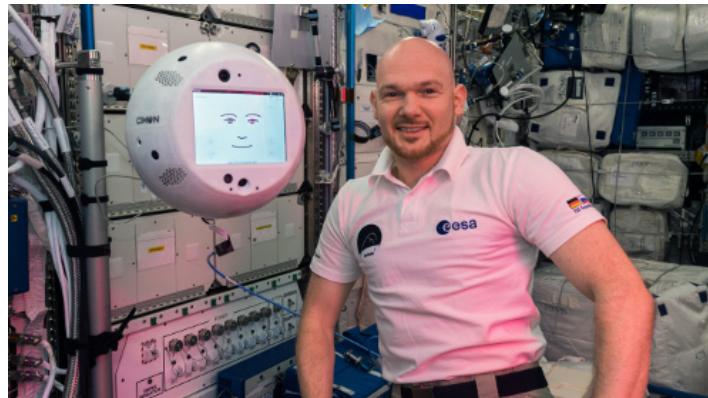
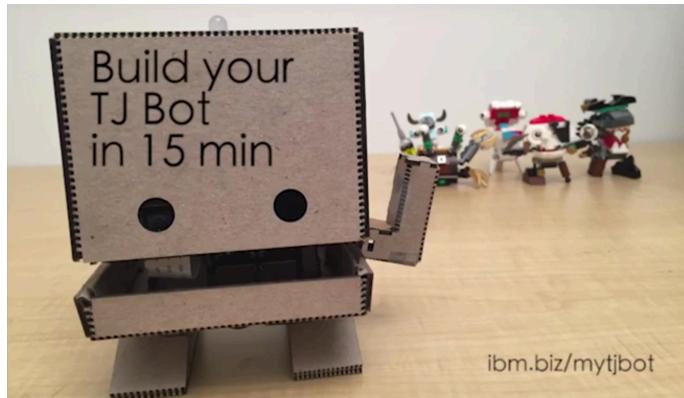
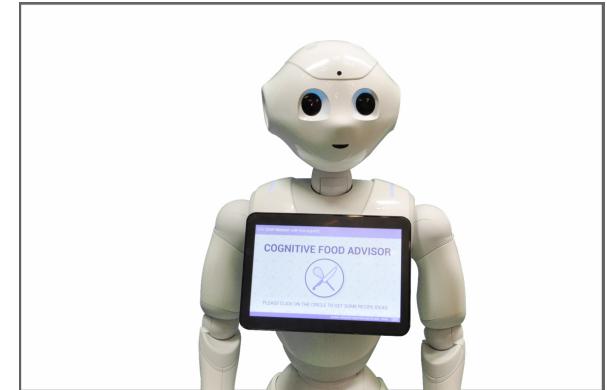
Model Agnostic Contrastive Explanations for Machine Learning Classification Models, Amit Dhurandhar et al  
<https://www.ibm.com/downloads/cas/0ZRZNR8E>

# Explainability



# Links and related Material

- IBM Watson OpenScale, Getting Started  
<https://cloud.ibm.com/docs/services/ai-openscale-icp?topic=ai-openscale-icp-gs-get-started>
- IBM AI Fairness 360 Framework (open source)  
<https://aif360.mybluemix.net/>
- IBM Watson Studio  
<https://cloud.ibm.com/catalog/services/watson-studio>
- IBM Cloud  
<https://cloud.ibm.com>



IBM

