

# Putting Artificial Intelligence Back into People's Hands

Toward an Accessible, Transparent and Fair AI

2 February 2020 · FOSDEM, Brussels, Belgium

Vincent Lequertier · FSFE Volunteer · <https://vl8r.eu>



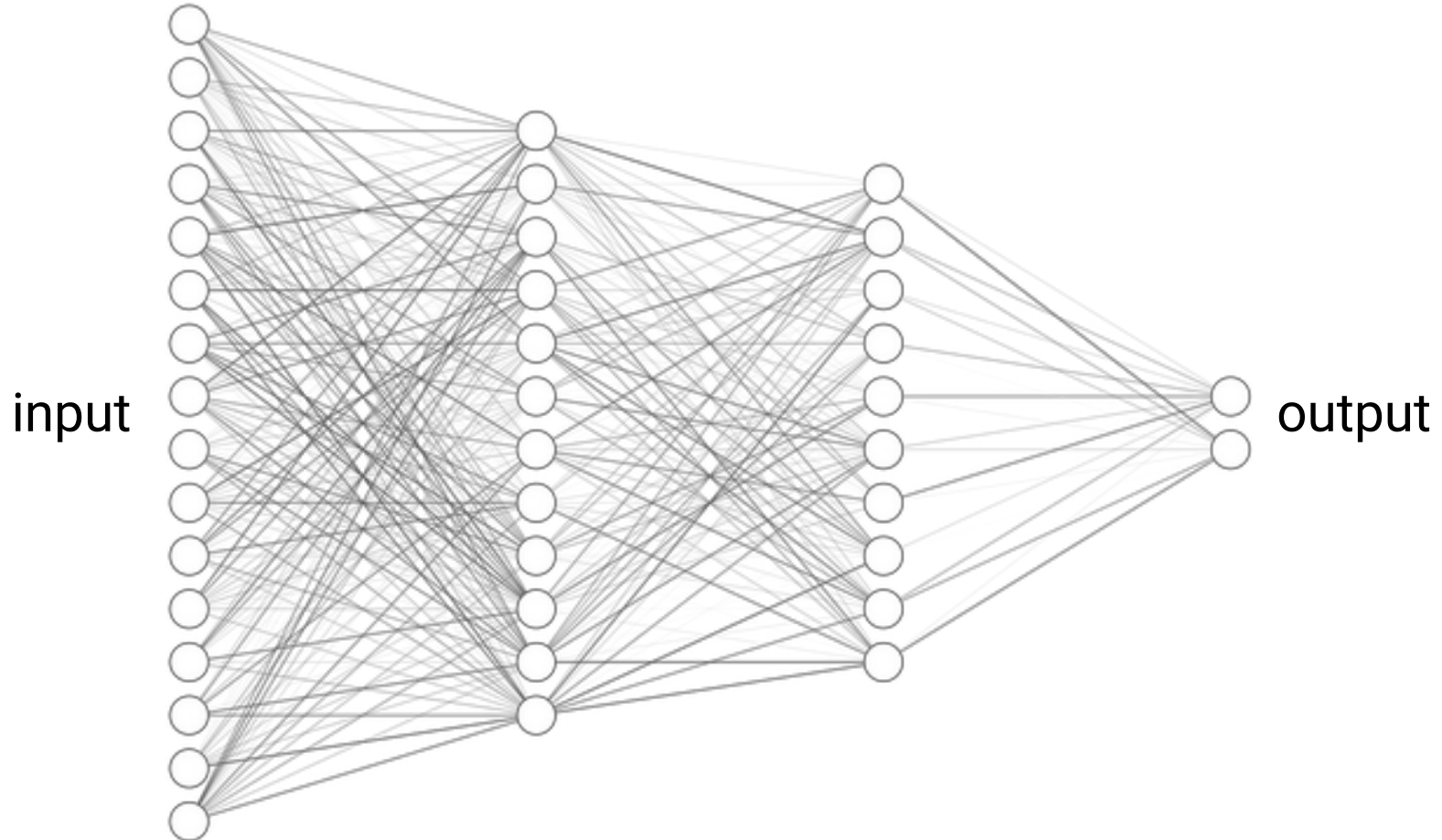
# Agenda

- How to create accessible Artificial Intelligence?
- Can AI be transparent and accurate?
- How to build fairness into AI?

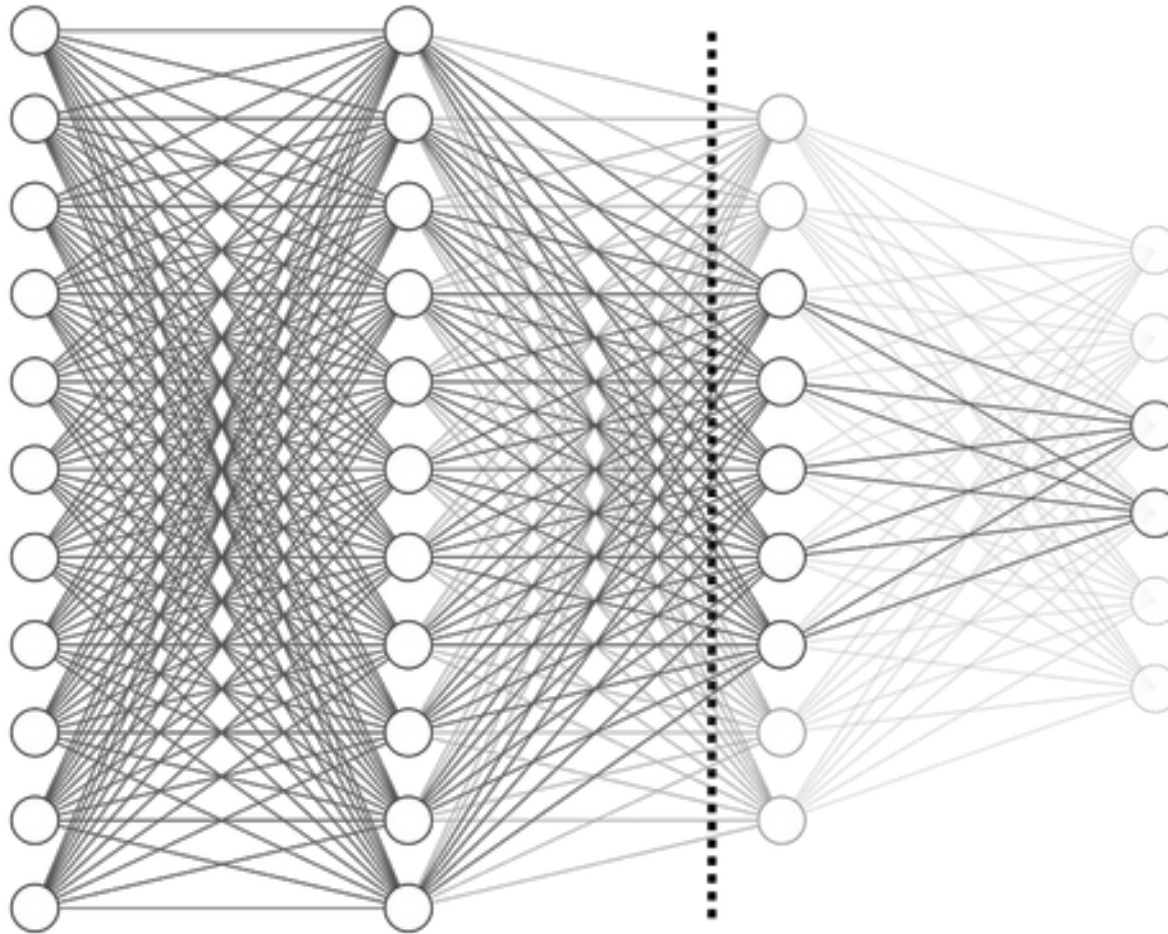


# Artificial Intelligence accessibility

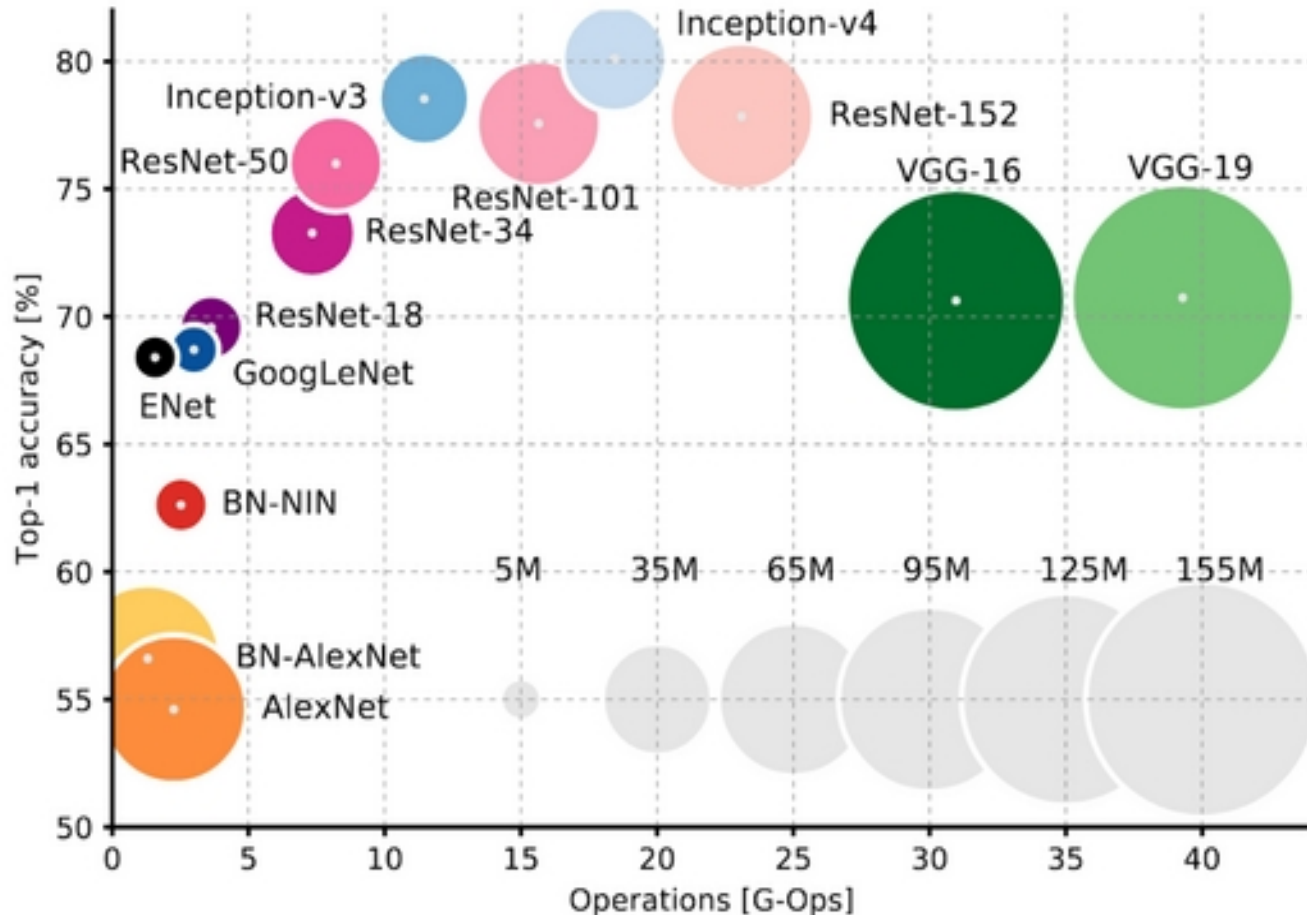
# What is a neural network?



# Leveraging other models: fine-tuning



# Bigger models are not more accurate



# How to make AI accessible?

- Make it easy to reuse the model (ONNX format)
- Release the training code and the dataset under a Free licence
- Consider the number of FLOP when designing the model



# Artificial Intelligence transparency



# AI is used for critical matters

- Loan approval
- Justice
- Healthcare
- Self-driving cars

# Why do we want transparency?

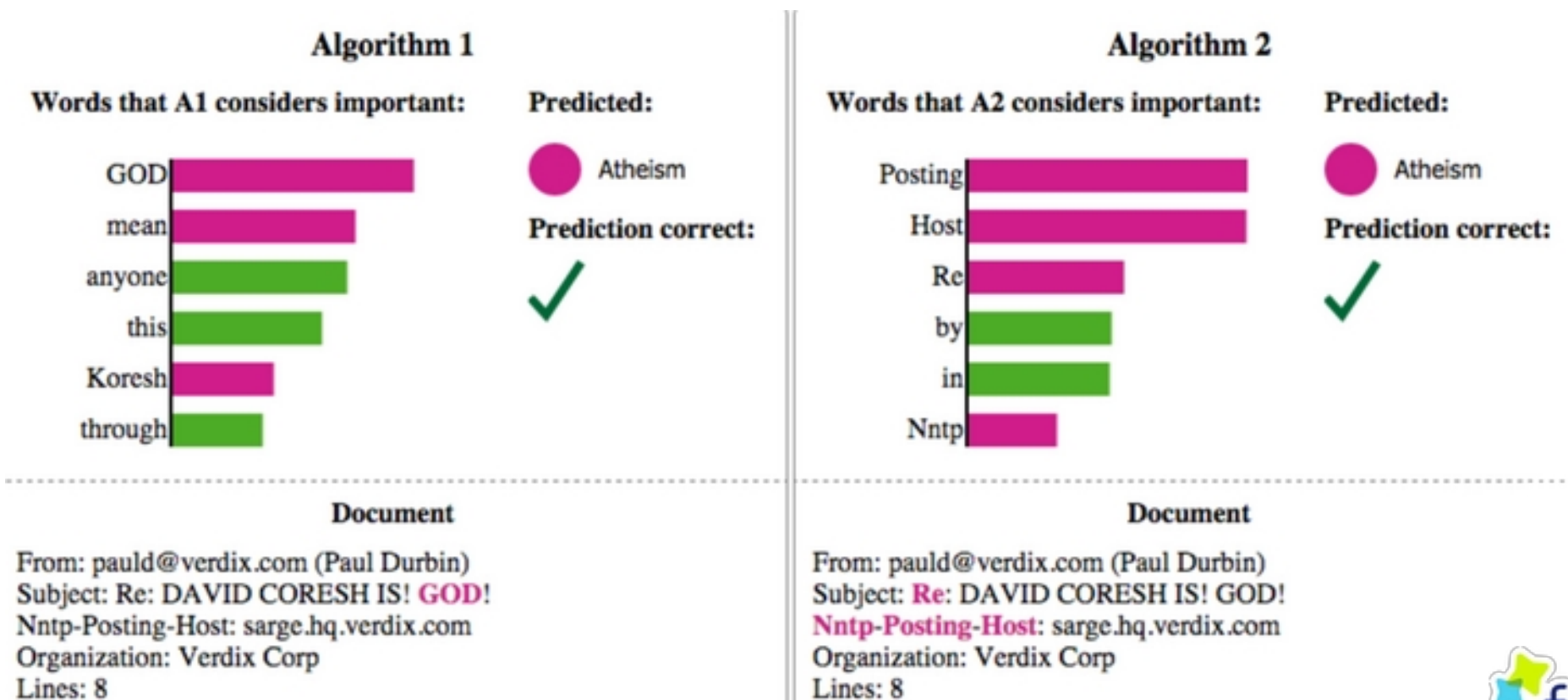
- Allows to interpret the result
- Builds trust in the model
- Makes debugging easier

# Parameters are not meant to be transparent



# LIME: Debugging and selecting models

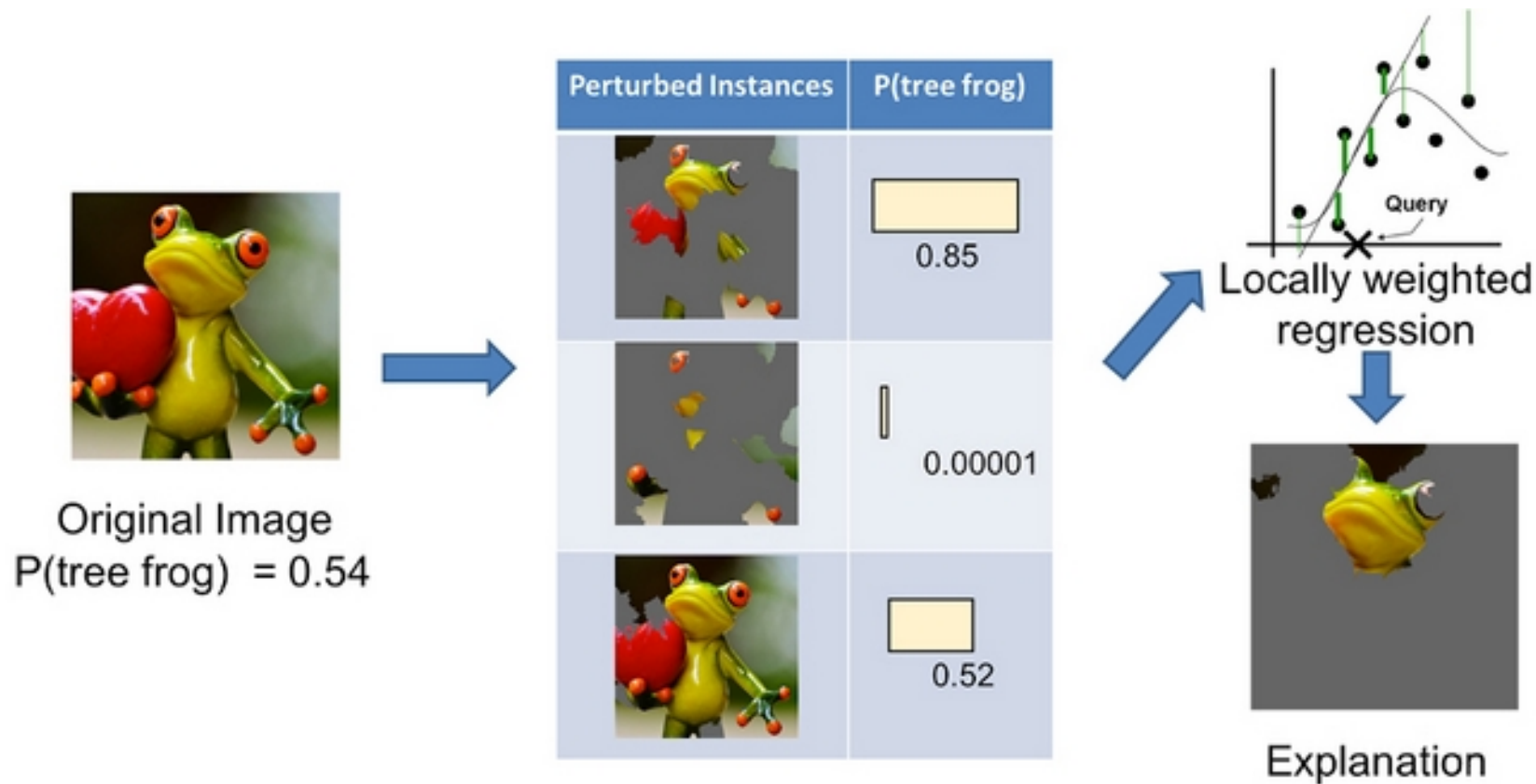
Local Interpretable Model-agnostic Explanations



# Making sense of images classification



# How does it work?

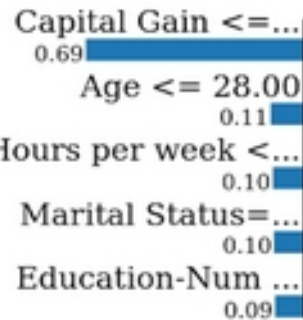


# Also for tabular data

Prediction probabilities



<=50K



>50K

Feature Value

Capital Gain	0.00
Age	19.00
Hours per week	30.00
Marital Status=Never-married	True
Education-Num	9.00



# Artificial Intelligence fairness



# Protecting car colors is easy

<b>brand</b>	<b>seats</b>	<b>year</b>	<b>color</b>	<b>speed (km/h)</b>
A	5	2011	blue	150
B	2	2012	black	200
C	5	2010	red	250

# Protecting gender is not easy

gender	hobby	education	salary
female	women's volleyball team	CS degree	35k
male	football team captain	self-taught	37k
male	chess	CS degree	37k

**⚠ Think about correlation before removing an attribute**

# Vocabulary

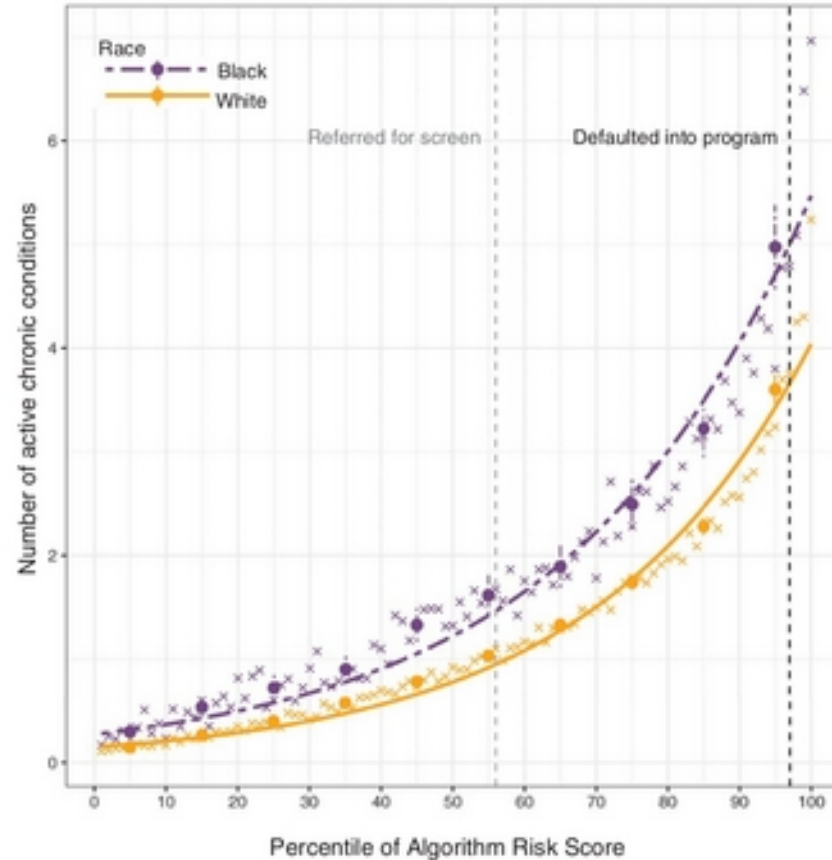
- True Positive (TP)
- True Negative (TN)
- False Positive (FP)
- False Negative (FN)

# COMPAS recidivism scoring

	All defendants		Black defendants		White defendants			
	Low	High	Low	High	Low	High		
Survived	2681	1282	Survived	990	805	Survived	1139	349
Recidivated	1216	2035	Recidivated	532	1369	Recidivated	461	505
<b>FP rate</b>	32.35		<b>FP rate</b>	44.85		<b>FP rate</b>	23.45	
<b>FN rate</b>	37.40		<b>FN rate</b>	27.99		<b>FN rate</b>	47.72	

propublica.org (2016)

# Racial bias in healthcare



Obermeyer, Z., Powers, B., Vogeli, C., & Mullainathan, S. (2019). Dissecting racial bias in an algorithm used to manage the health of populations.

# Why an algorithm can be unfair?

- Bias in the input data itself
- Training with the wrong metric (bias by proxy)
- Bad prediction model
- Bias is hard to notice
- "*With great power comes great responsibility*" (Peter Parker)

# A fair loss function

Let  $k$  be the number of values of a protected attribute

Let  $f : y_{pred}, y_{true} \rightarrow s \in [0, 1]$  be a fairness function

$$loss = loss + \lambda \frac{\sum_{i=0}^k w_i f_i(y_{pred}, y_{true})}{\min_{\forall i \in [0, k[} f_i(y_{pred}, y_{true})}$$

# Thank you! Questions?

2 February 2020 · FOSDEM, Brussels, Belgium

Vincent Lequertier · FSFE Volunteer · <https://vl8r.eu>

