

Explainable AI (XAI): Confronting Bias, Discrimination, and Fairness in Machine Learning

Michael Ridley

Postgraduate Affiliate, Vector Institute
PhD Candidate, FIMS, Western University
Librarian Emeritus, University of Guelph

mridley@uoguelph.ca @mridley

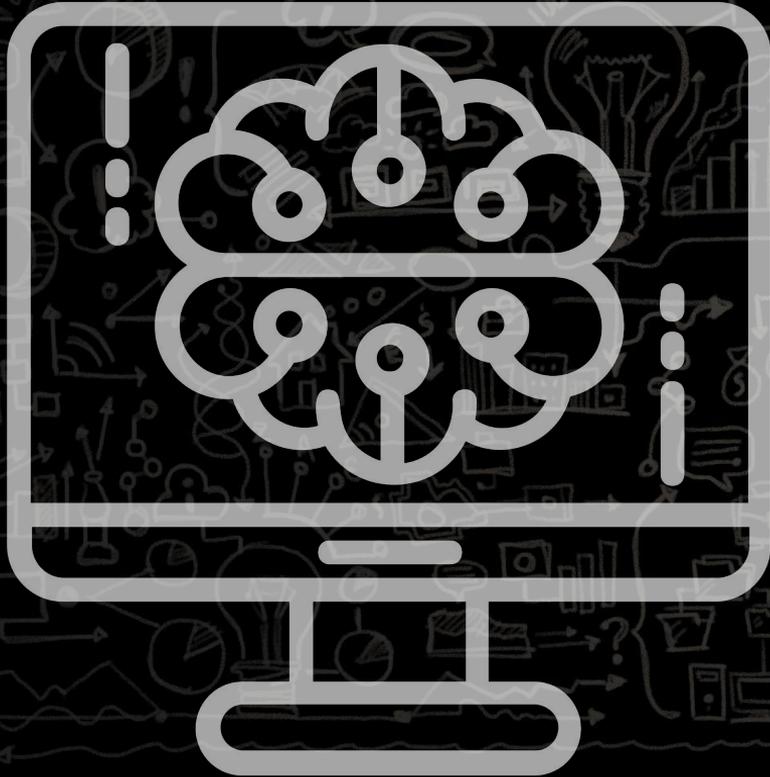
Access 2019 – University of Alberta – October 2019



“I believe that artificial intelligence will become a major human rights issue in the twenty-first century.”

Safiya Noble,
Algorithms of Oppression (2018)

Machine Learning



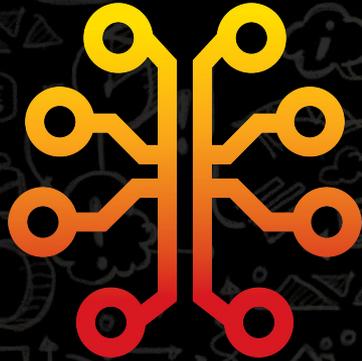
Ubiquitous

Powerful

Opaque

Invisible

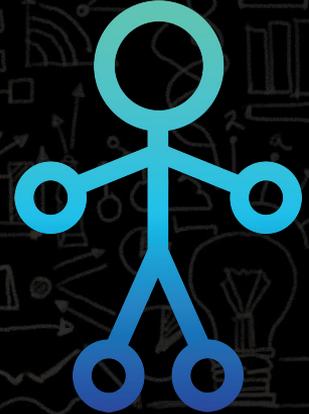
Consequential



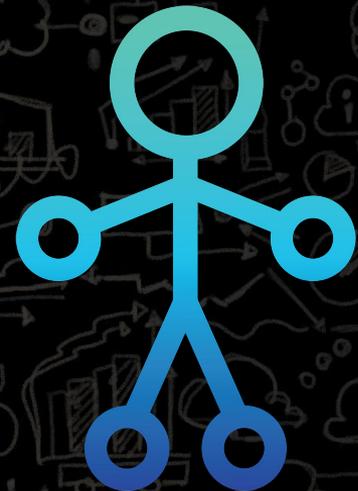
“The danger is not so much
in delegating cognitive tasks,

but in distancing ourselves from
– or in not knowing about –

the nature and precise mechanisms
of that delegation”



de Mul & van den Berg (2011). Remote control: Human
autonomy in the age of computer-mediated agency.



What do people want
in an explanation?

How can AI provide such an explanation?

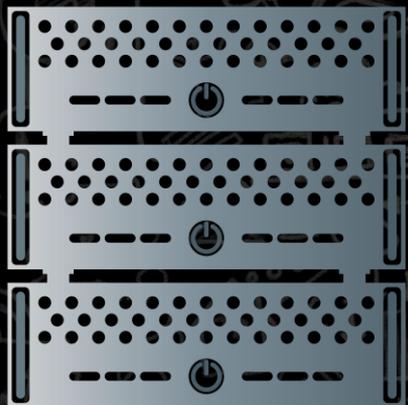


General Data
Protection Regulation
(GDPR)

The image features a dark blue background with a subtle gradient. In the center, there is a circle of twelve yellow five-pointed stars, arranged in a ring. The text "right to explanation" is written in white, sans-serif font across the middle of the stars. The text is enclosed in quotation marks, with the opening quote on the left and the closing quote on the right, both positioned between the two stars that are horizontally aligned with the text.

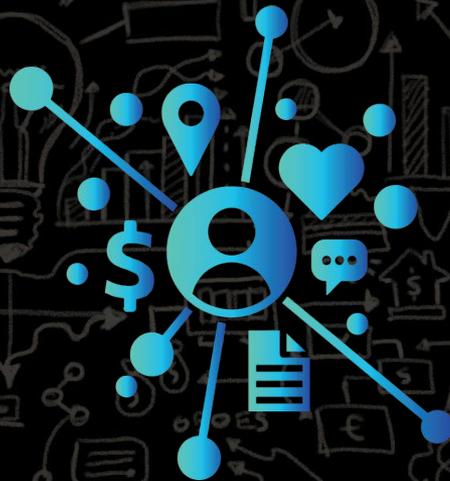
“right to explanation”





Computation

+



Data

+



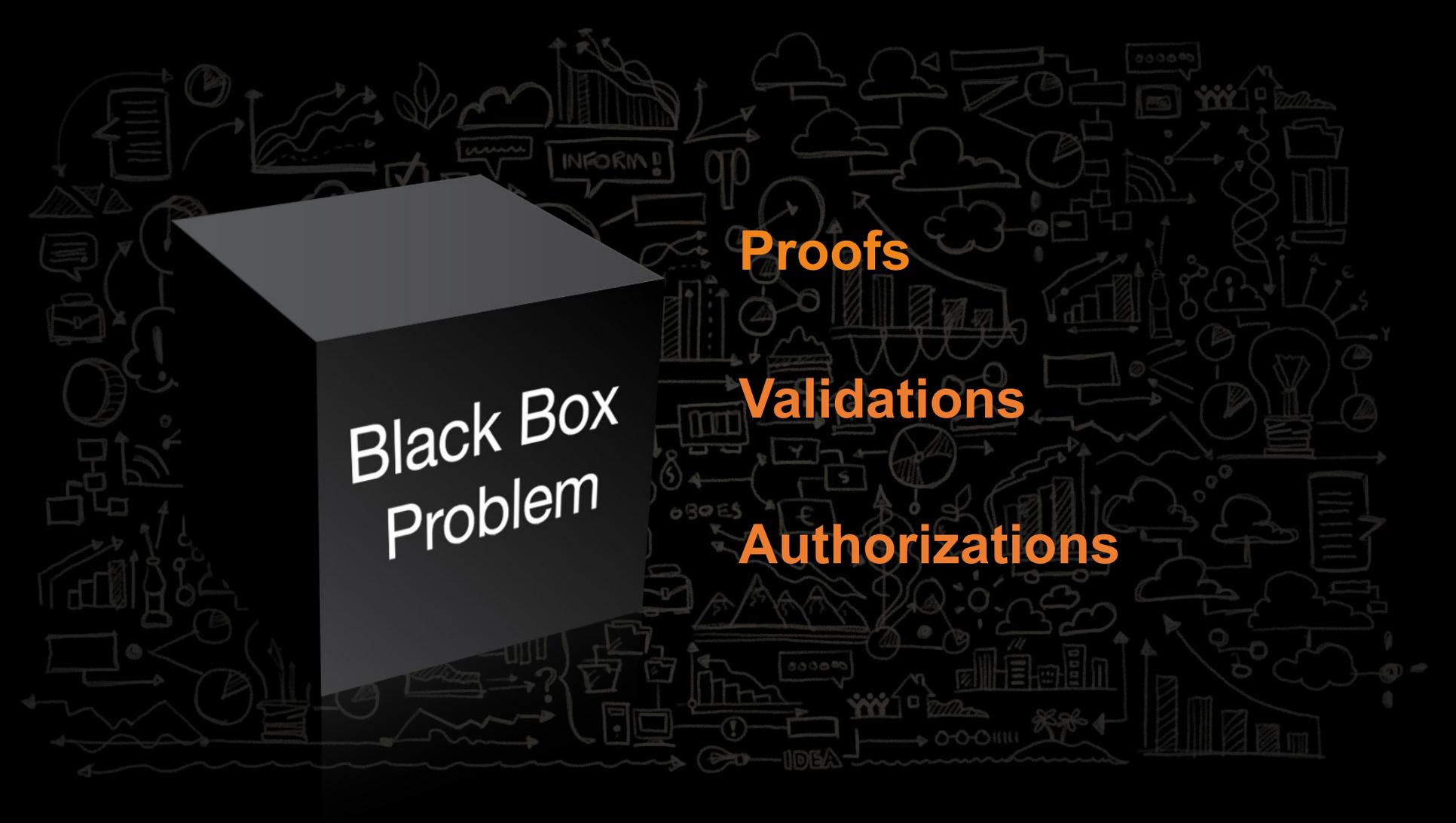
Algorithms

Social

Political

Economic



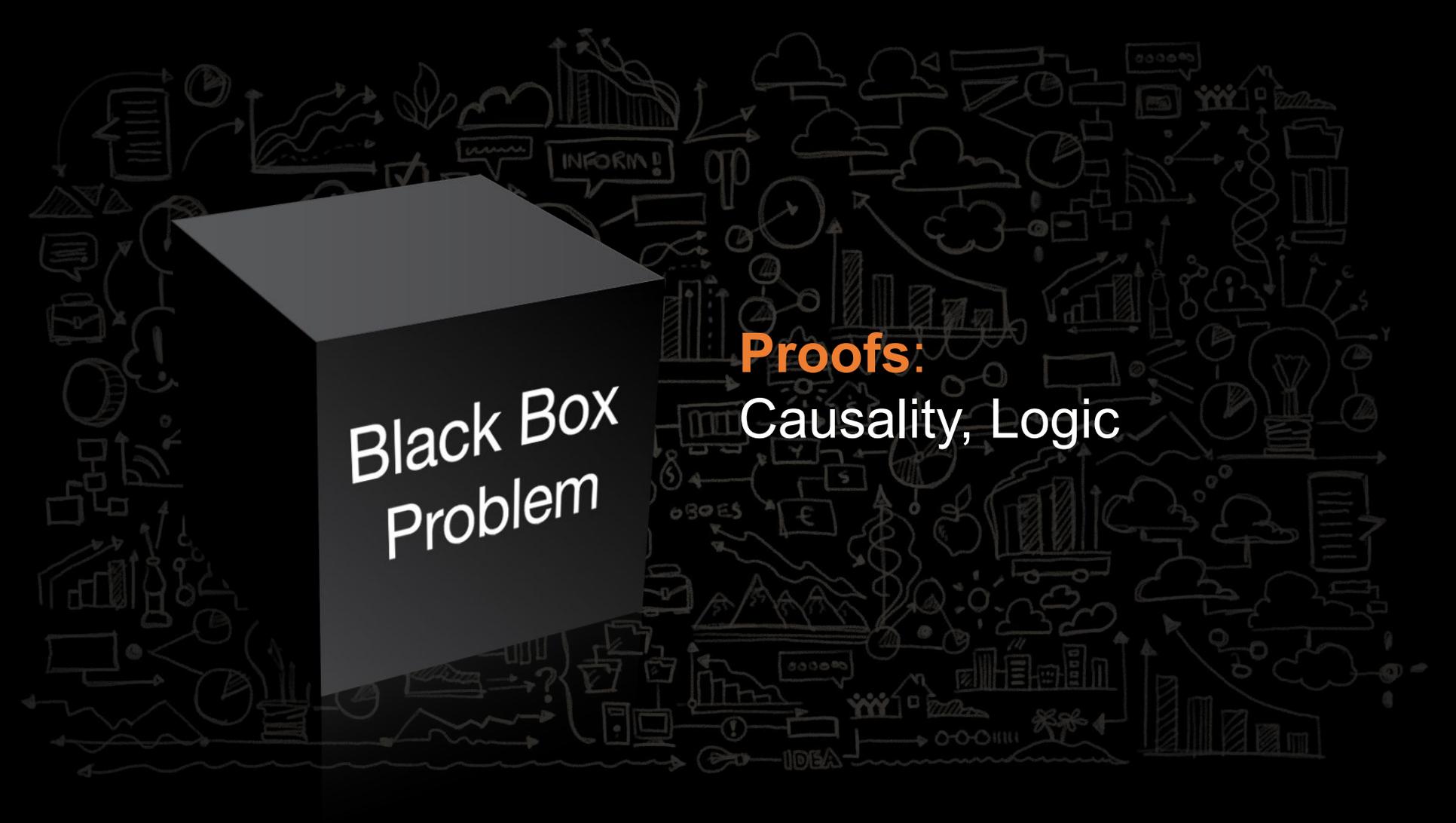


Black Box
Problem

Proofs

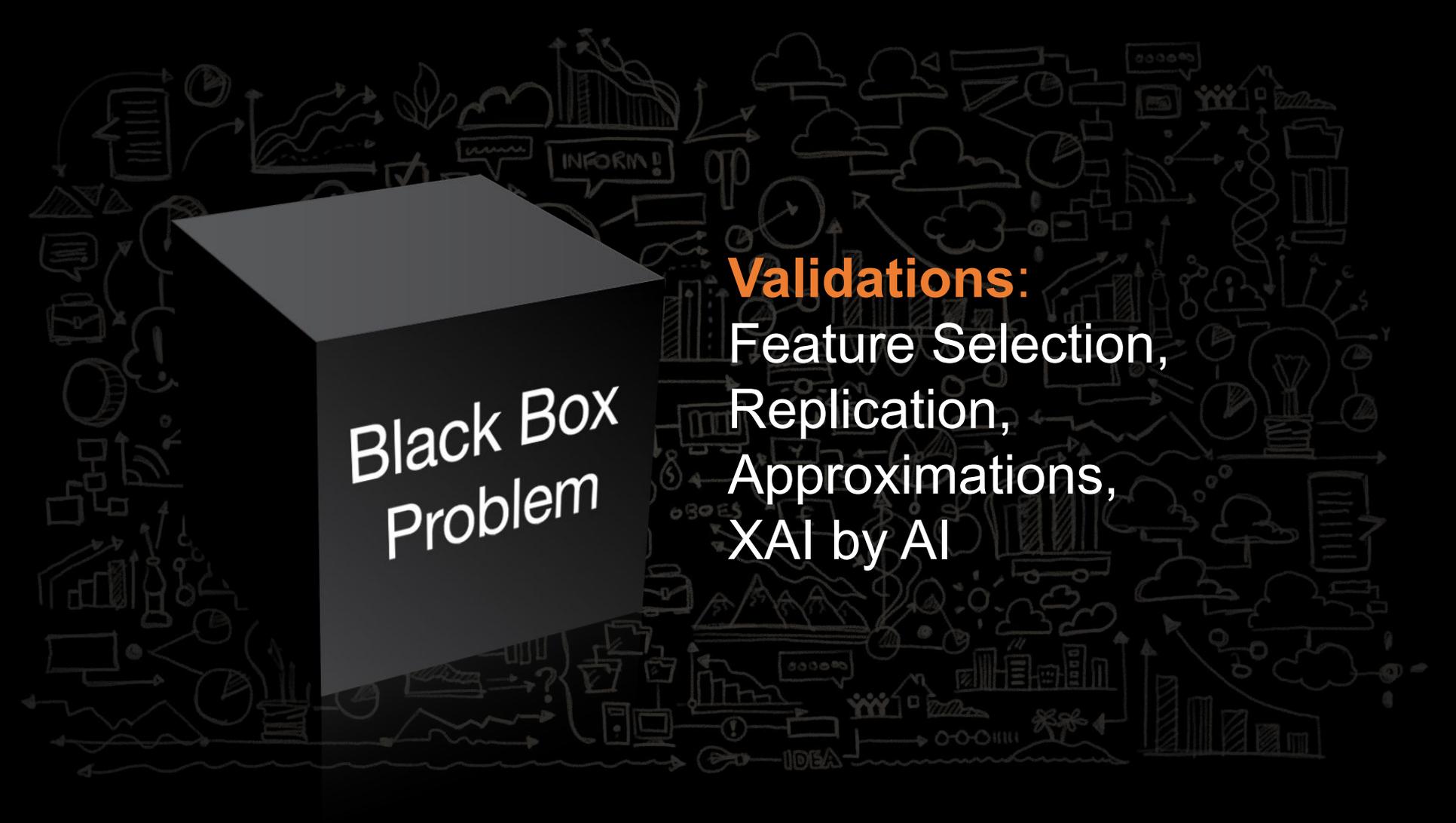
Validations

Authorizations



Black Box
Problem

Proofs:
Causality, Logic



Black Box
Problem

Validations:

Feature Selection,
Replication,
Approximations,
XAI by AI



Black Box
Problem

Authorizations:

Codes & Standards,
Legislation, Due Process,
Audit, Regulation



Roles for Libraries

Algorithmic Literacy

ML Partnerships & Development

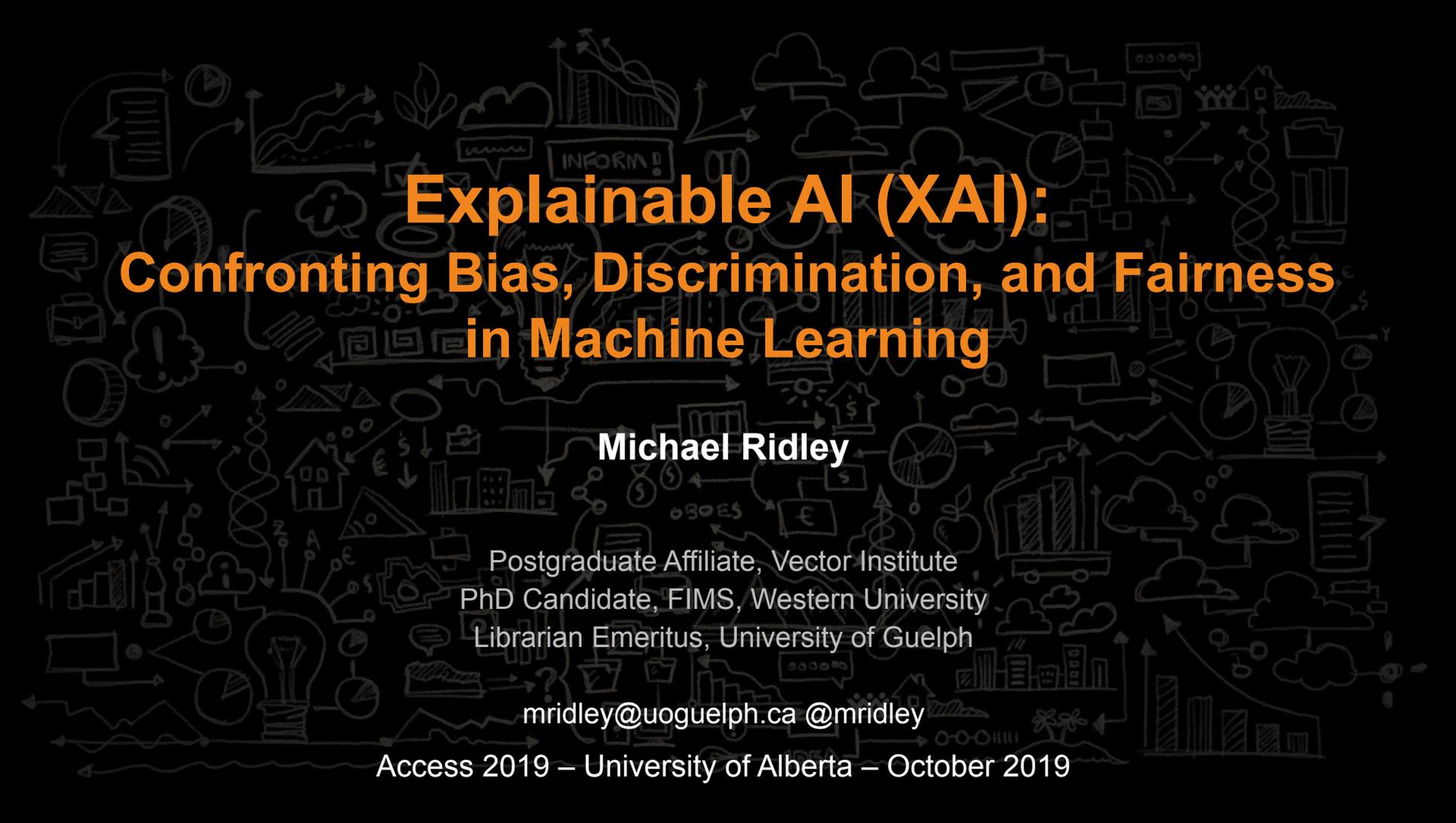
Advocacy

The New Digital Divide

“A class of people
who can use algorithms
and a class
used by algorithms”

David Lankes, Director
School of Library and Information Science
University of South Carolina





Explainable AI (XAI): Confronting Bias, Discrimination, and Fairness in Machine Learning

Michael Ridley

Postgraduate Affiliate, Vector Institute
PhD Candidate, FIMS, Western University
Librarian Emeritus, University of Guelph

mridley@uoguelph.ca @mridley

Access 2019 – University of Alberta – October 2019