# Cyber System Emulation for Dataset Development

Jamie Thorpe, jthorpe@sandia.gov

*Cybersecurity R&D, Sandia National Laboratories*
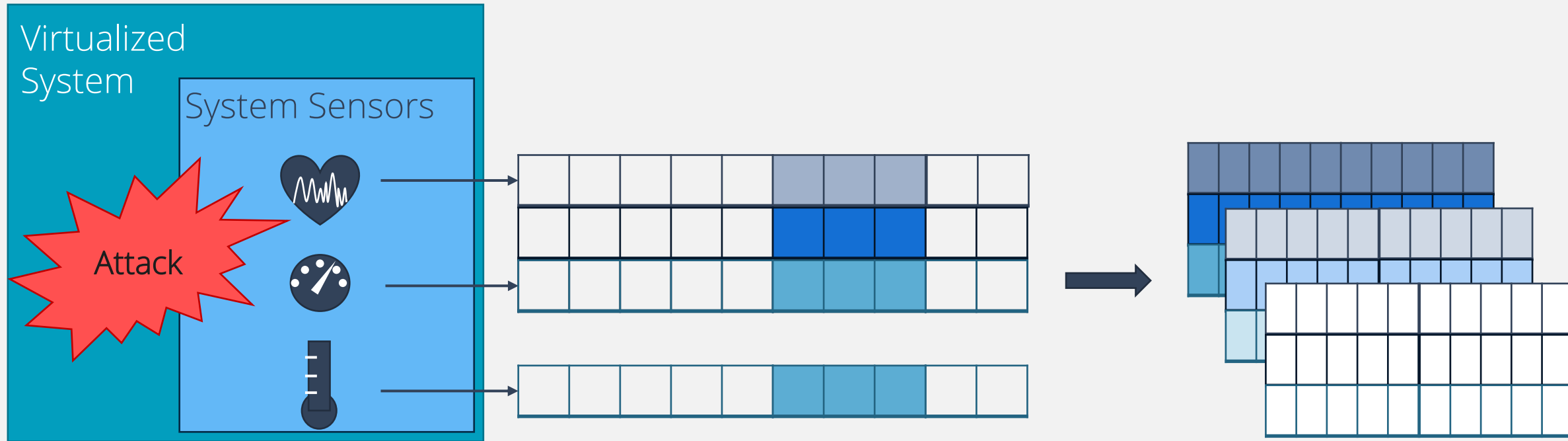
International Conference on Machine Learning

Machine Learning for Cybersecurity Workshop

Baltimore, MD, USA

July 22, 2022

# What is "Emulation"?



ACTUAL SYSTEM

VIRTUALIZED TESTBED

SIMULATION

"BAD DAY" BRAINSTORMING

REAL HARDWARE
REAL SOFTWARE

ABSTRACT HARDWARE
REAL SOFTWARE

ABSTRACT HARDWARE
ABSTRACT SOFTWARE

SUBJECT MATTER
EXPERT-DRIVEN

# Emulation and Data Collection

# Emulation for Dataset Development

What are some of the big challenges for system emulation?

- Emulation model creation and validation

- Emulation verification

- Data collection

  - What data can be collected?

  - How much data do we collect?

  - How do we know that the data collection process doesn't disrupt the emulation?

# Emulation for Dataset Development

Why isn't emulation the "simple" solution to developing datasets?

- Real-time runs

- Can Generative ML be applied here? If so, we need to consider...
    - How to apply common generative algorithms to multivariate timeseries
        - Different structure
        - Different datatypes
        - Different data relationships
    - How many samples do we need from the emulation in order to accurately represent each class?

# There's a lot to think about!

Emulation isn't THE solution, but it could be PART of the solution.