# PANEL 6: METRICS AND MEASUREMENT METHODS: WHAT AND HOW TO TEST

**José Hernández-Orallo** (jorallo@upv.es)
Valencian Research Institute for Artificial Intelligence (vrAIn) (vrain.upv.es)
Universitat Politècnica de València, València (www.upv.es)
Leverhulme Centre for the Future of Intelligence, Cambridge (lcfi.ac.uk)

*NIST Workshop on AI Measurement and Evaluation, June 15-17, 2021*

**NIST Special Publication 970**

# Measuring the Performance and Intelligence of Systems: Proceedings of the 2000 PerMIS Workshop
## August 14-16, 2000

Edited by:
A. M. Meystel
E. R. Messina
*Intelligent Systems Division*
*Manufacturing Engineering Laboratory*
*National Institute of Standards and Technology*
*Gaithersburg, MD 20899-8230*

Co-Sponsored by:
National Institute of Standards and Technology
Defense Advanced Research Projects Agency
Institute of Electrical and Electronic Engineers Control Systems Society
National Aeronautics and Space Administration

In Cooperation with:
Institute of Electrical and Electronic Engineers Neural Network Council

HERNÁNDEZ-ORALLO

THE MEASURE OF ALL MINDS

# THE MEASURE OF ALL MINDS

*Prose Award 2018*

**Evaluating Natural and Artificial Intelligence**

CAMBRIDGE
UNIVERSITY PRESS
www.cambridge.org

ISBN 978-1-107-15301-1

CAMBRIDGE

**JOSÉ HERNÁNDEZ-ORALLO**

# AI Evaluation as Aggregated Performance

- **GOAL:** Estimate the expected result $\tilde{R}$ of system $\pi$ and a new task $\mu$.

Given:

- Distribution $p$ in problem class $M$ (e.g., configurations of a navigation task)
- Metric of performance $\mathbb{R}$ (e.g., navigation success)

**Calculate** aggregated performance and **extrapolate** for $\mu$!

$$\tilde{R}(\boldsymbol{\pi}, \boldsymbol{\mu}) = \sum_{\boldsymbol{\mu}' \in \mathrm{M}} p(\boldsymbol{\mu}') \mathbb{R}(\boldsymbol{\pi}, \boldsymbol{\mu}')$$

- This is useful **if $\mu \sim p$ and** the operating condition in $\mathbb{R}$ does not change.
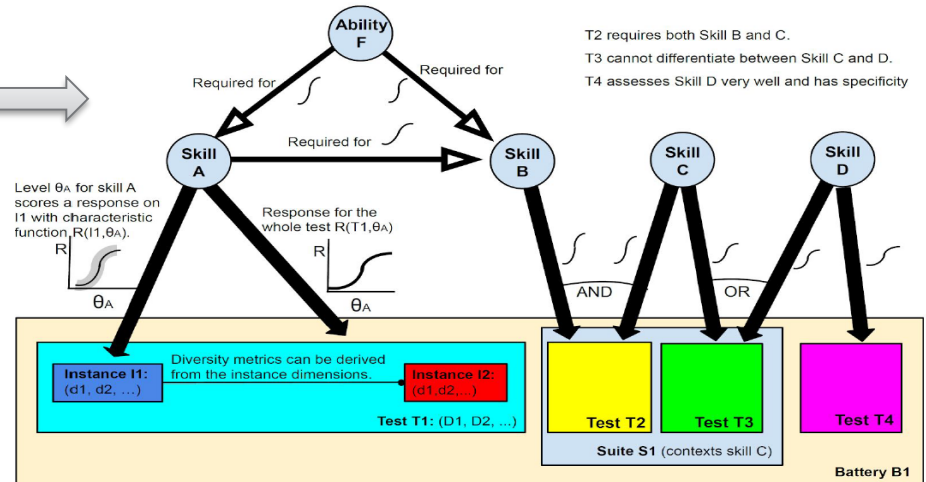
But this is almost never the case!

- Infers a **capability profile** for system $\pi$. We also use a **problem profile** for $\mu$.

$\pi$

Introductory
Internal Models
Y-Mazes
Weak Generalisation
Detour Tasks
Tool Use
Radial Mazes
Numerosity
Spatial Elimination
Object Permanence
Delayed Gratification
Support & Gravity

$\mu$

Item Characteristic Curves

- Given both, estimate $\tilde{R}(\boldsymbol{\pi}, \boldsymbol{\mu})$
- Key ideas:
  - Instance difficulties become dual to capabilities (à la IRT).
  - Requires identifying the capabilities and their relation.
  - Constructs for $\pi$ and $\mu$ are latent factors: measurement is no longer additive.

# ReCOG-AI : Measurement Layouts

- **Robust Evaluation of Cognitive Capabilities and Generality in AI**
  - 2021-2023 (planning to work with DARPA)
    - related to the machine common sense program, director: Matt Turek.
  - Run at the Centre for the Future of Intelligence, Cambridge, UK.
  - Measurement Layouts:
    -
  - Generality:
    - In RL settings for basic navigation skills
    - With language or multimodal models

# ReCOG-AI : Spaces and Features

**128x128 RGB pixels**



- **Original feature space:**
  - observable by the system. Usually abstracted into latent features.

- **Surface feature space:**
  - sometimes observable. A general system should be invariant to these.

**Symmetry**     **Irrelevant elements**



- **Cognitive (construct) space:**
  - usually non-observable. Performance should correlate with them:
    - agents with a high capabilities profile in this space will imply success for problems with lower difficulty levels in these capabilities.

# METRICS AND MEASUREMENT

- **Metrics**:
  - Capabilities should have a proper scale.
  - Aggregations are not additive from results.
    - More detailed results, annotated instances!
    - No more aggregated results only, please!
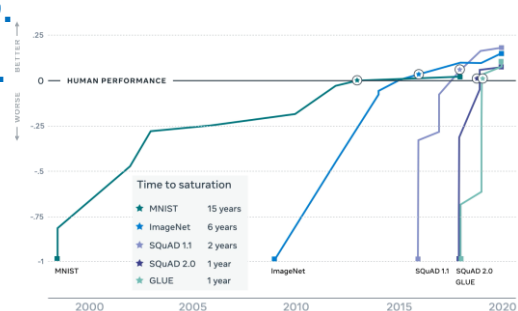    - No more "superhuman" claims, please!
- **Measurement**:
  - Cover the capabilities space, not the original *p*.
  - Avoid "challenge-solve-and-replace" dynamics.
  - Explore instance variation:
    - Adaptive testing
    - Adversarial testing

Hernandez-Orallo, J. "AI Evaluation: On Broken Yardsticks and Measurement Scales", MetaEval@AAAI2020.



ImageNet competition test set accuracy — ImageNet 2012 validation set accuracy — Human performance

ImageNet competition ends in 2017.



AI benchmark saturation over time

CIFAR10 → CIFAR100,
SQuAD1.1 → SQuAD2.0,
GLUE → SUPERGLUE,
Starcraft → Starcraft II

"Give me the data (distribution) and I will ace the test in a year!"

Time to saturation
- MNIST       15 years
- ImageNet    6 years
- SQuAD 1.1   2 years
- SQuAD 2.0   1 year
- GLUE        1 year

THANKS!

# OTHER SOURCES AND INITIATIVES:

- Other Talks (http://josephorallo.webs.upv.es/)
  - Diversity Unites Intelligence: Measuring Generality
  - Measuring A(G)I Right: Some Theoretical and Practical Considerations
  - Natural and Artificial Intelligence: Measures, Maps and Taxonomies
- Book (http://allminds.org):
  - The Measure of All Minds: Evaluating Natural and Artificial Intelligence, Cambridge University Press 2017
- The AI Collaboratory: http://aicollaboratory.org/
  - Part of the European Commission's AI watch:
    - https://ec.europa.eu/knowledge4policy/ai-watch_en
- ReCOG-AI and the animal AI environment:
  - Part of the Kinds of Intelligence Programme at the CFI in Cambridge
    - http://lcfi.ac.uk/projects/kinds-of-intelligence