

Capability-Oriented Evaluation (of AI)

Jose Hernandez-Orallo

Mostly based on

Ryan Burnell, John Burden, Danaja Rutar, Konstantinos Voudouris, Lucy Cheke, Jose Hernandez-Orallo
“Not a Number: Identifying Instance Features for Capability-Oriented Evaluation“ IJCAI 2022



BEYOND THE IMITATION GAME: QUANTIFYING AND EXTRAPOLATING THE CAPABILITIES OF LANGUAGE MODELS

performance

Do we have capability-oriented evaluation in AI?

Alphabetic author list:

Aarohi Srivastava, Abhinav Rastogi, Abhishek Rao, Abu Awal Md Shoeb, Abubakar Abid, Adam Fisch, Adam R. Brown, Adam Santoro, Aditya Gupta, Adria Garriga-Alonso, Agnieszka Klaska, Aitor Lewkowycz, Akshat Agarwal, Alethea Power, Alex Ray, Alex Warstadt, Alexander W. Kocurek, Ali Safiya, Ali Tazarar, Alice Xiang, Alicia Parrish, Allen Nie, Aman Hassain, Amanda Askell, Amanda Dsouza, Ambrose Stone, Ameet Rahane, Anantharaman S. Iyer, Anders Andreassen, Andrea Madotto, Annemarie Santilli, Andreas Stuhlmüller, Andrew Dai, Andrew La, Andrew Lampinen, Andy Zou, Angela Jiang, Angelica Chen, Anh Vuong, Animesh Gupta, Anna Gottardi, Antonio Norelli, Ana Venkatesh, Arash Gholami-avoodi, Arfa Tabassum, Arul Menezes, Arun Kirubakaran, Asher Mullikandov, Ashish Vasudeva, Austin Herrick, Avia Efrat, Aykut Erdem, Ayla Karakaş, B. Ryan Roberts, Bao Sheng Loc, Barret Zoph, Bartłomiej Bojanowski, Batuhan Özyurt, Behnam Hedayatnia, Behnam Neyshabur, Benjamin Inden, Benno Stein, Berk Ekmecki, Bill Yuchen Lin, Blake Howald, Cameron Diaz, Cameron Dour, Catherine Sinson, Cedrick Arqueta, César Ferri Ramirez, Chandan Singh, Charles Rathkopf, Chenlin Meng, Chitta Baral, Chiyu Wu, Chris Callison-Burch, Chris Waites, Christian Voigt, Christopher D. Manning, Christopher Potts, Cindy Ramirez, Clara E. Rivera, Clemencia Siro, Colin Raffel, Courtney Ashcraft, Cristina Garbacea, Damien Sileo, Dan Garrette, Dan Hendrycks, Dan Kilman, Dan Roth, Daniela Freeman, Daniel Khoshabi, Daniel Levy, Daniel Mosegui González, Danielle Perszyk, Danny Hernandez, Danqi Chen, Daphne Ippolito, Dar Gilboa, David Dohan, David Drakard, David Jurgens, Debajyoti Datta, Deep Ganguli, Denis Emelin, Denis Kleyko, Deniz Yuret, Derek Chen, Derek Tam, Dheevuk Hupkes, Diganta Misra, Dilyara Duzan, Dimitris Cocchio Molli, Diyi Yang, Dong Ho Lee, Ekaterina Shustova, Ekin Dogus Cubuk, Elad Nofal, Eleanor Hagerman, Elizabeth Barnes, Elizabeth Donoway, Ella Pavlick, Emanuele Rodola, Emma Lam, Eric Chu, Eric Tang, Erkut Erdem, Ernie Chang, Ethan A. Chi, Ethan Dyer, Ethan Jerzak, Ethan Kim, Eunice Engeru Manyasi, Evgenii Zheltonozhskii, Fanyue Xia, Fatemeh Siar, Fernando Martínez-Plumed, Francesca Happé, François Chollet, Frieda Rong, Gaurav Mishra, Genta Indra Winata, Gerard de Melo, Germán Kruszewski, Giambattista Parascandolo, Giorgio Mariani, Gloria Wang, Gonzalo Jaimovich-López, Gregor Betz, Guy Gur-Ari, Hanu Galijasevic, Hannah Kim, Hannah Rashkin, Hannah Hajishirzi, Harsh Mehta, Hayden Bogar, Henry Shveth, Hinrich Schütze, Hiroo Yikura, Hongming Zhang, Hugh Mee Wong, Ian Ng, Isaac Noble, Jaap Jumelet, Jack Geisinger, Jackson Kernion, Jacob Hilton, Jaehoon Lee, Jaime Fernández Fisac, James B. Simon, James Koppel, James Zhang, James Zou, Jan Kocoň, Jana Thompson, Jared Kaplan, Jarema Radom, Jascha Sohl-Dickstein, Jason Phang, Jason Wei, Jason Yosinski, Jekaterina Novikova, Jelle Bosscher, Jennifer Marsh, Jeremy Kim, Jeron Tal, Jesse Engel, Jessjoba Alabi, Jiacheng Xu, Jiaming Song, Jillian Tang, Joan Waweru, John Burden, John Miller, John J. Bennis, Jonathan Berant, Jörg Froberg, Jos Rozen, Jose Hernandez-Orallo, Joseph Boudeman, Joseph Jones, Joshua B. Tenenbaum, Joshua S. Rule, Joyce Chua, Kim Kanceler, Karen Livescu, Karf Krauth, Karthik Gopalakrishnan, Katerina Ignatyeva, Katja Markert, Kaustubh D. Dhole, Kevin Gimpel, Kevin Omondo, Kory Mathewson, Kristen Chafiallo, Ksenia Shkrikina, Kumar Shridhar, Kyle McDonell, Kyle Richardson, Lara Reynolds, Le Gao, Li Zhang, Liam Dugan, Lianhui Qin, Lidia Conterato, Louis-Philippe Morency, Luca Moschella, Lucas Lam, Luca Noble, Ludwig Schmidl, Luheng He, Luis Oliveros Colón, Luke Metz, Lütfi Kerem Senel, Maarten Bosma, Maarten Sap, Maarje ter Hoeve, Maheen Farooqi, Manaal Faruqi, Mantas Mazeika, Marco Baturan, Marco Marelli, Marco Maru, Maria Jose Ramirez Quintana, Marie Tokikhe, Mario Giudianelli, Martha Lewis, Martin Pothast, Matthew L. Leavitt, Matthias Hagen, Mátyás Schubert, Medina Orduna Baitemirova, Mody Anand, Melvin McElrath, Michael A. Yee, Michael Cohen, Michael Gu, Michael Ivanitskiy, Michael Sturritt, Michael Strube, Michal Śwadowski, Michele Bevilacqua, Michihito Yasunaga, Mihir Kale, Mike Cain, Milice Xu, Mírac Suzgun, Mo Tiwari, Mohit Bansal, Moira Aminasari, Mox Geva, Mozdeh Ghahini, Mukund Varma T, Nanyun Peng, Nathan Chi, Nayoon Lee, Neta Gur-Ari Krakover, Nicholas Cameron, Nicholas Roberts, Nick Douron, Nikita Nangia, Niklas Deckers, Niklas Muennighoff, Nitish Shirish Keskar, Niveditha S. Iyer, Noah Constant, Noah Fiedel, Nuan Wen, Oliver Zhang, Omar Agha, Omar Elbaghdadi, Omer Levy, Owain Evans, Pablo Antonio Moreno Casares, Parth Doshi, Pascale Fung, Paul Pu Liang, Paul Vicol, Pegah Alipoormolabashi, Peiyuan Liao, Percy Liang, Peter Chang, Peter Eckersley, Phu Mon Htut, Pinyu Hwang, Piotr Miłkowski, Prayush Patil, Pooja Pezeshkpour, Priti Oli, Quozhu Mei, Qing Lyu, Qianlong Chen, Rabin Banjale, Rachel Etta Rudolph, Raefel Gabriel, Rahul Habacker, Ramón Risco Delgado, Raphaël Millière, Rhythm Garg, Richard Barnes, Rif A. Saurous, Riku Arakawa, Robbe Raymaekers, Robert Frank, Rohan Sikand, Roman Novak, Roman Sitelew, Ronan LeBras, Rosanne Liu, Rowan Jacobs, Rui Zhang, Ruslan Salakhutdinov, Ryan Chi, Ryan Lee, Ryan Stovall, Ryan Teehan, Ryan Yang, Sahib Singh, Saif M. Mohammad, Sajant Anand, Sam Dillavou, Sam Shleifer, Sam Wiseman, Samuel Gruetter, Samuel R. Bowman, Samuel S. Schoenholz, Sanghyun Han, Sanjeev Kwatra, Sarah A. Rous, Sarik Ghazarian, Sayan Ghosh, Sean Casey, Sebastian Bischoff, Sebastian Gehrmann, Sebastian Schuster, Sepideh Sadeghi, Shadi Hamdan, Sharon Zhou, Shashank Srivastava, Sherry Shi, Shikhar Singh, Shima Asaadi, Shiyang Shene Gu, Shubh Pachchigar, Shubham Toshnival, Shayan Uppalhy, Shyamolika Shanmugam (Shammie), Debanshu Siamak Shakeri, Simon Thomeyer, Simone Melzi, Siva Reddy, Sneha Priscilla Makini, Soo-Hwan Lee, Spencer Torene, Sriharsha Hatwar, Stanislas Dehaene, Stefan Dvíc, Stefane Ermon, Stella Biderman, Stephanie Lin, Stephen Prasad, Steven T. Piantadosi, Stuart M. Shieber, Summer Mishnerghi, Svetlana Kiritchenko, Swaroop Mishra, Tal Linzen, Tal Schuster, Tao Li, Tao Yu, Tariq Ali, Tatsuo Hashimoto, Te-Lin Wu, Théo Desbordes, Theodore Rothchild, Thomas Pan, Tianle Wang, Tiberius Nkinyili, Timo Schick, Timofei Kornev, Timothy Telleen-Lawton, Titus Tundany, Tobias Gerstenberg, Trenton Chang, Trishay Neeraj, Tushar Khot, Tyler Shultz, Uri Shaham, Vedant Misra, Vera Demberg, Victoria Nyamai, Vikas Raunak, Vinay Ramasesh, Vinay Uday Prabhu, Vishakh Padmakumar, Vivek Srivastava, William Fedus, William Saunders, William Zhang, Wout Vossen, Xiang Ren, Xiaoyu Tong, Xinran Zhao, Xinyi Wu, Xudong Shen, Yadollah Yaghoobzadeh, Yair Lektrez, Yanguo Song, Yasmine Bahri, Yejin Choi, Yichi Yang, Yiding Hao, Yifu Chen, Yonatan Belinkov, Yu Hou, Yufang Hou, Yuntao Bai, Zachary J. Wu, Zhaoye Zhao, Zijian Wang, Ziji J. Wang, Ziru Wang, Ziyi Wu

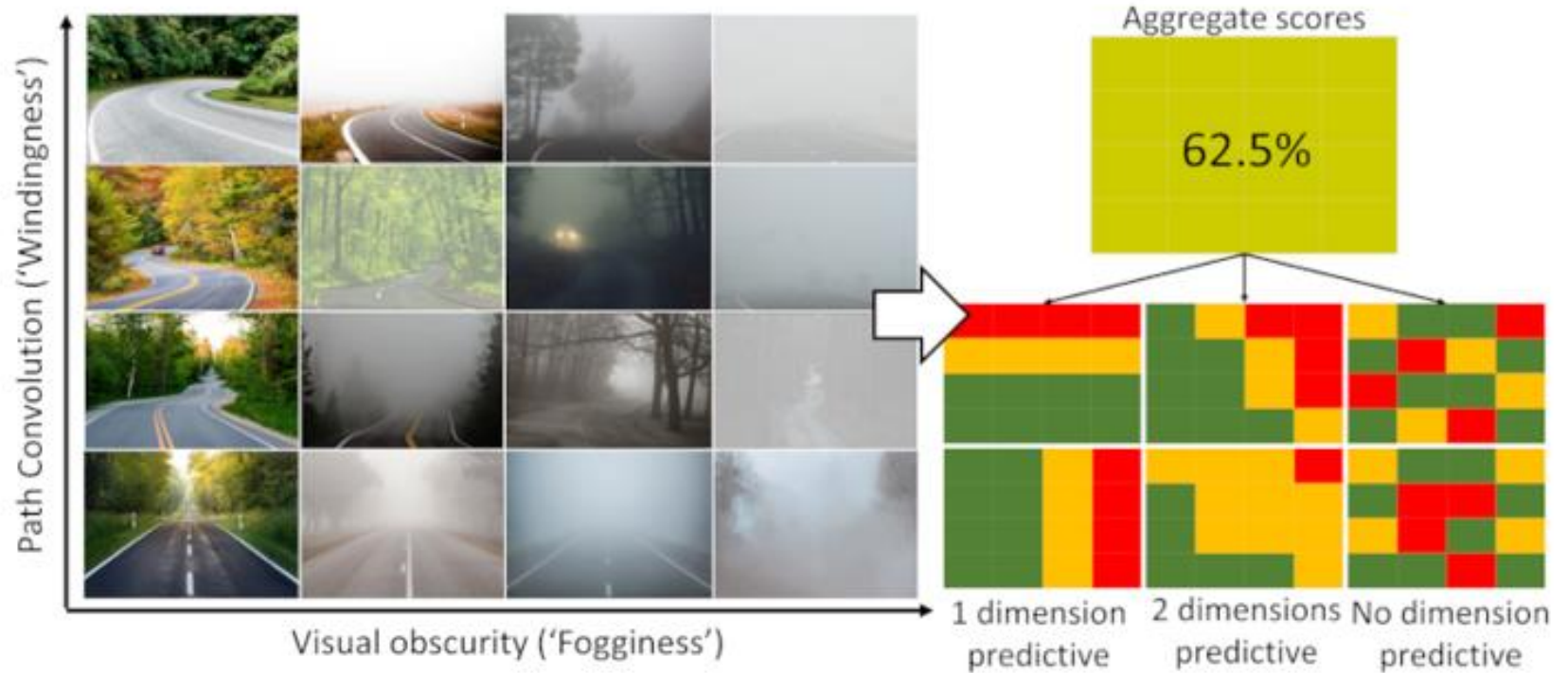
Performance-oriented vs Capability-oriented

- Performance is a property (a measure) of a **pair <system, item>**:
 - Examples:
 - Correct prediction of MySpamFilter on Email735
 - 85% accuracy of ResNet23 on ImageNet
 - **Performance changes when the item/distribution changes**
 - On blurry, adversarial, OOD images the result is much worse
- Capability is a **property of a system**:
 - Examples:
 - The system can add integers up to three digits.
 - The system can jump up to 1.20 metres high.
 - **Capability doesn't change when the item/distribution changes**
 - Bar at 1.50 metres high? Bad performance because the capability is lower.

The problems of aggregated performance

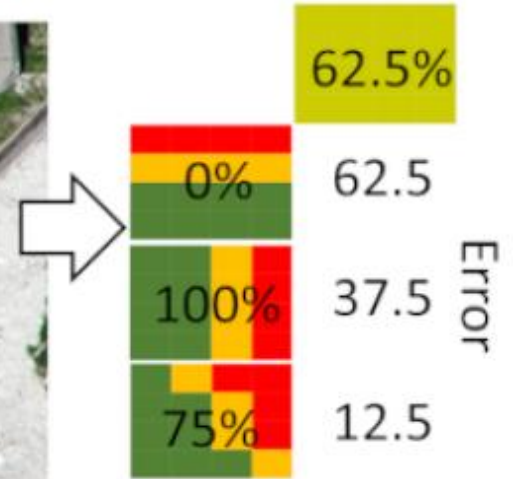
- No patterns of performance
 - No identification of failure points
- Poor estimation of performance for new distributions
 - The metric cannot be extrapolated
- Poor granular estimation for the same distribution!
 - Likely to be conditions under which the system performs better or worse

Where does my system fail?



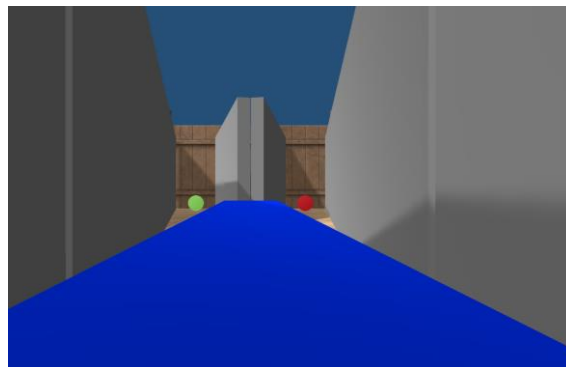
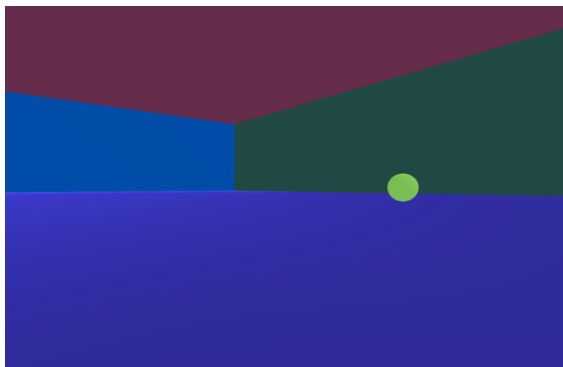
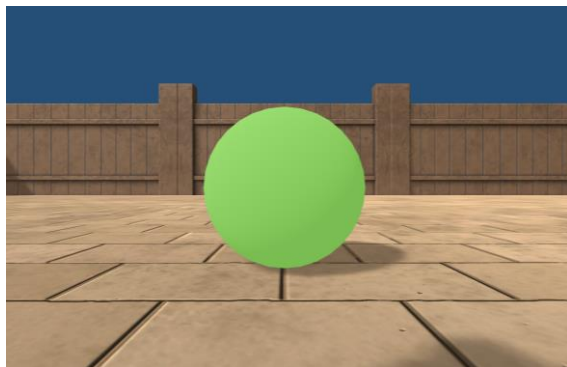
Estimate success/failure granularly

Can I safely drive this road on a clear day?



Proof of concept: Animal AI Olympics

- Selected subset of AAI/O instances measuring simple goal-directed behaviour
- Data across 99 instances from 68 agents



<http://animalaiolympics.com/AAI/>

M Crosby, B Beyret, M Shanahan, J Hernández-Orallo, L Cheke, M Halina “The animal-AI testbed and competition” NeurIPS 2019 Competition and Demonstration Track, Proceedings of Machine Learning Research, 2020

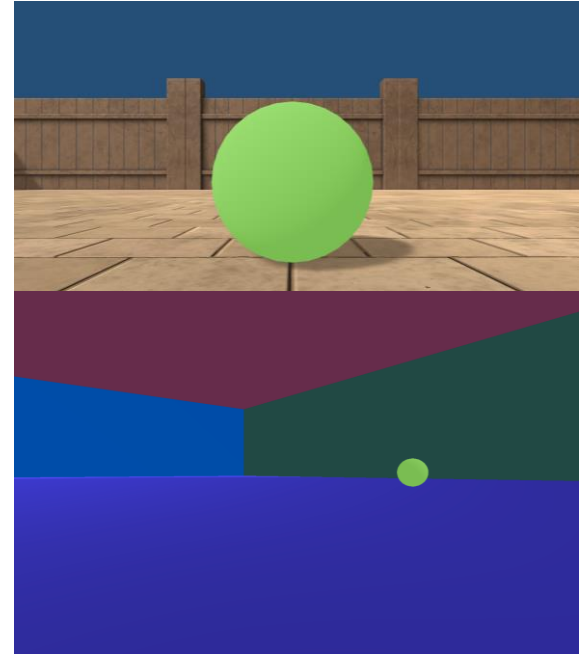
Identifying features of interest

Relevant

- Reward size
- Reward distance
- Reward in view (i.e., in front vs behind)

Irrelevant

- Reward side (left vs right)
- Reward colour (green vs yellow)



Identified dimensions and agent characteristic curves

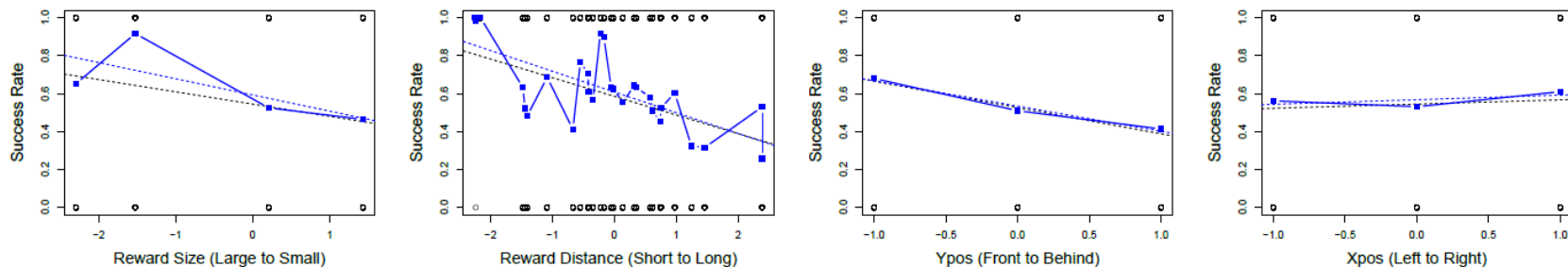


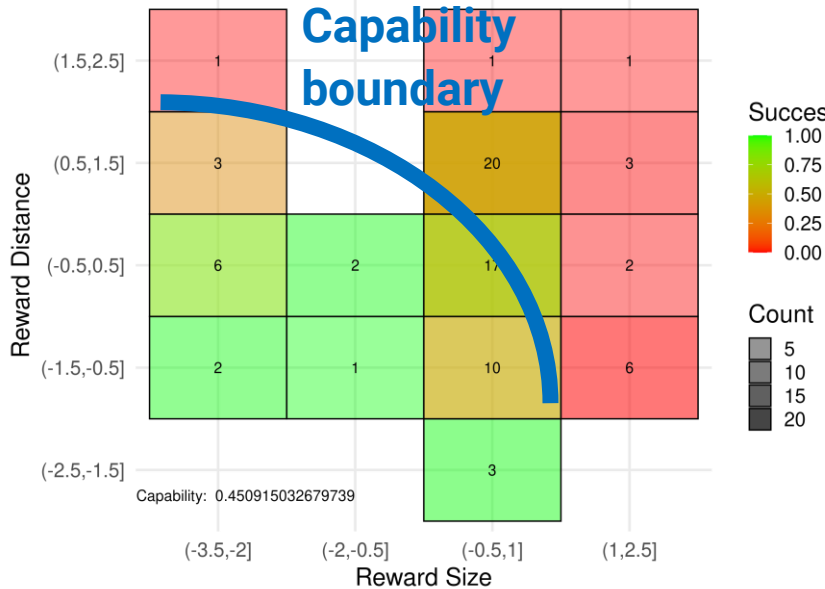
Figure 5: Characteristic curves of all competition entrants (agents) according to three relevant features (size, distance and Ypos) and one irrelevant feature (Xpos). Black dashed lines show the linear regression for the black points (pass/fail), while blue dashed lines interpolate the blue points (binned success rate). The conformances (Spearman correlations against monotonic sequence) are 0.80, 0.60, 1.00 and -0.50 , respectively.

Visualising performance distribution

- Plot a subset of relevant variables:

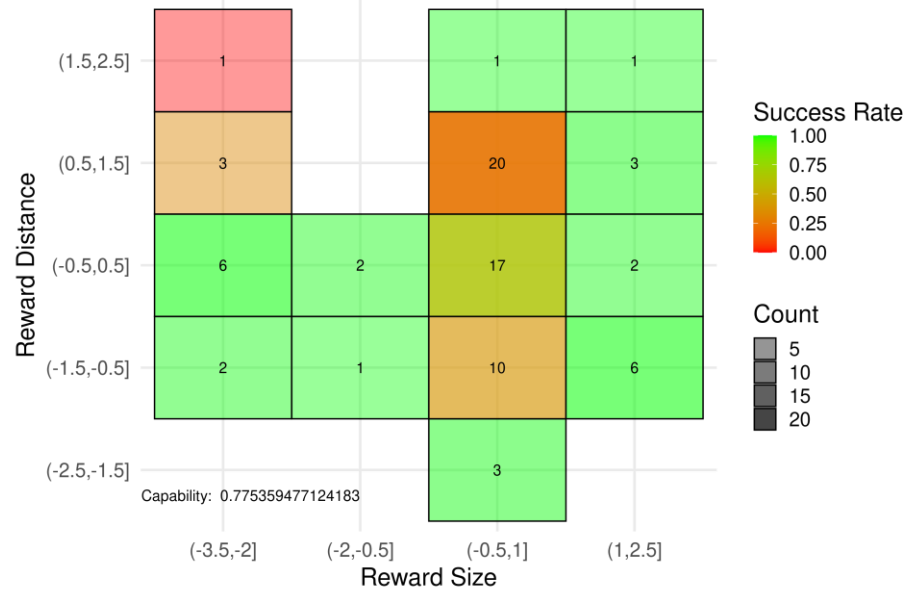
This system doesn't show monotonicity. We can't identify any level of capability robustly.

Juohmaru



Conformant System

y.yang



Non-Conformant System

Predicting performance

extrapolate Global
ACCuracy
? = 54.7%

(ignores system locality
and feature relevance)

extrapolate agent
ACCuracy
? = 46.8%

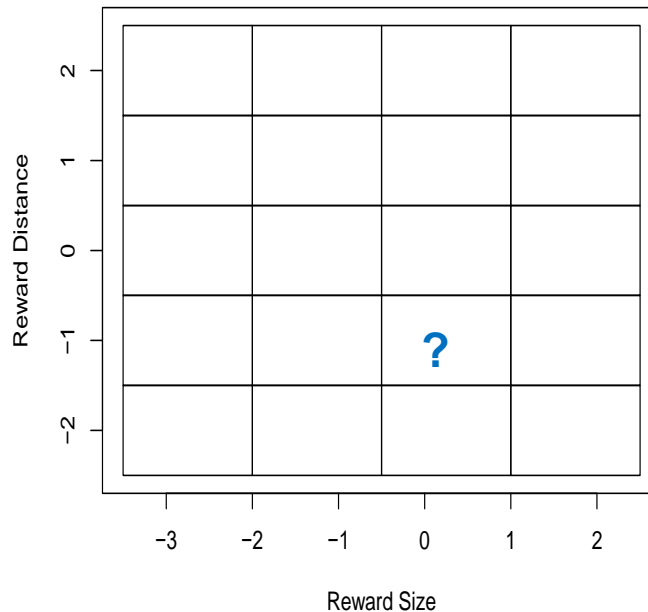
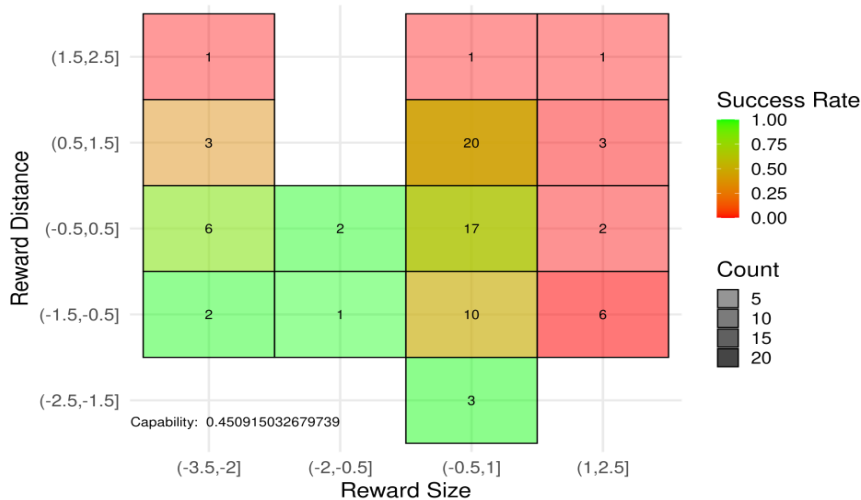
(ignores feature
relevance)

extrapolate bin
? = 40%

(ignores other bins)

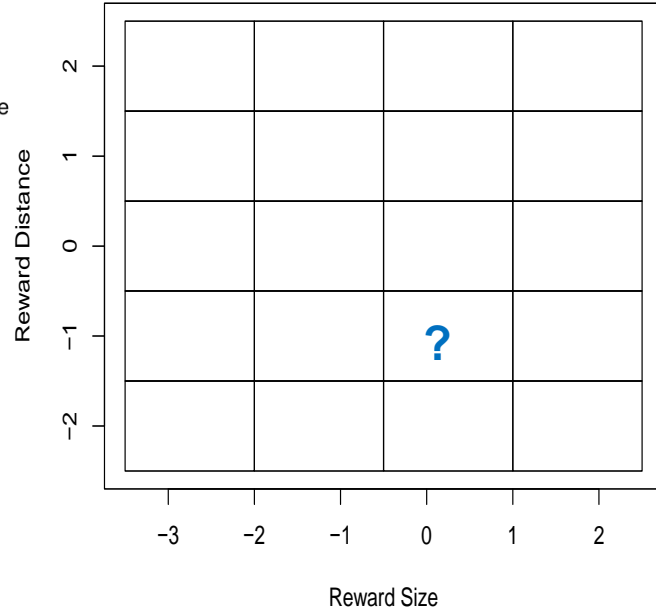
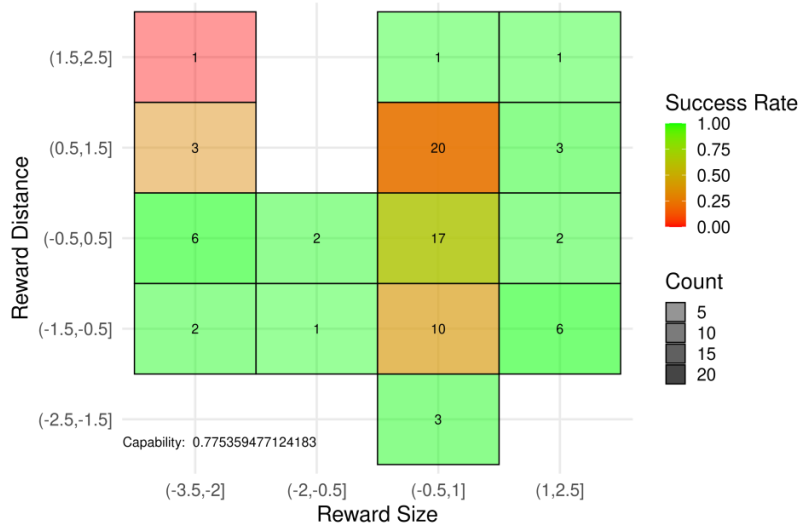
use parametric model
on capabilities
? = 73.2%

(parameter goodness-of-fit
may be poor)



Except the last one, these are basically non-inferential methods (constant models or binning extrapolations)

Predicting performance (parametric model won't fit)



use assessor
models
(Using all variables or
only the relevant ones?)

assessors = let's use all the power of ML to characterise the system's performance!!

Assessors (non-parametric models)


Hernández-Orallo, J.; Schellaert, W.; Martínez Plumed “Training on the Test Set: Mapping the System-Problem Space in AI”, AAAI 2022 (Blue Sky Ideas Award).

- Conditional probability estimator of the result r for AI system π on situation μ :

$$\hat{R}(r|\pi, \mu) \approx \Pr(R(\pi, \mu) = r)$$

- It is trained (and evaluated) on test data:
 - Using a distribution of situations (instances) μ .
 - Using a distribution of systems π .

It is applied during deployment, before π does any inference or even starts.



π	μ	r
Resnet, $\theta_1, \theta_2, \dots$	Image3, χ_1, χ_2, \dots	1
Resnet, $\theta_1, \theta_2, \dots$	Image23, χ_1, χ_2, \dots	0
...
Inception, $\theta_1, \theta_2, \dots$	Image3, χ_1, χ_2, \dots	1
Inception, $\theta_1, \theta_2, \dots$	Image78, χ_1, χ_2, \dots	1
...

Predicting performance (Comparison)

	Maj. (1)	G.Acc.	T.Acc.	~All+A	~Rel+A
Error	45.3%	48.0%	33.6%	19.7%	20.6%
MAE	45.3%	49.6%	34.9%	29.3%	30.2%
MSE	45.3%	24.8%	17.6%	14.8%	15.4%

Animal AI Competition Data: 99 instances x 68 agents

Guidelines for Capability-based Evaluation

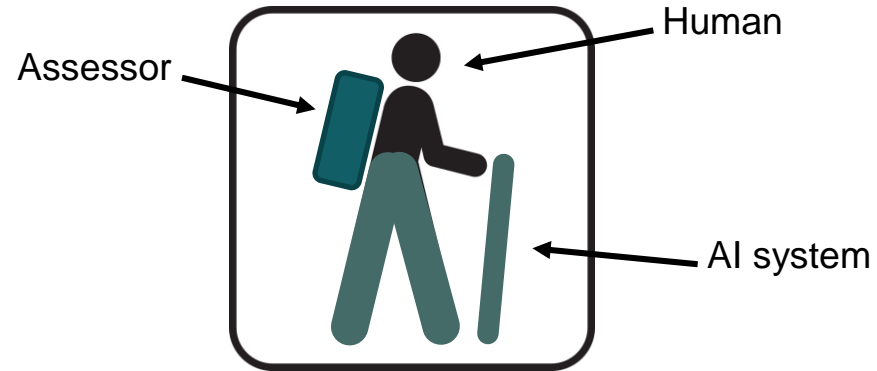
1. Choose a domain, task or benchmark with instance-level data.
2. Identify features that can be extracted or easily annotated for each instance during testing.
3. Identify which features are relevant (should affect performance) and those that are irrelevant (should not affect performance) based on theory and domain knowledge.
4. Analyse the relationships between features and performance using correlation and other exploratory analyses.
5. Select features of interest, bin them appropriately and build characteristic grids (both global and for individual systems) to evaluate patterns of performance.
6. Build predictive models using extracted features. Compare the results with other ways of predicting performance, such as extrapolating average metrics.
7. Using characteristic grids and predictive models, evaluate capabilities of each system across the distributions of the dimensions of interest.
8. Identify areas of competence for individual systems so that these can later be used for testing more complex or advanced capabilities and skills (where appropriate).
9. Use areas of weakness to inform changes to the benchmark, system models or training.
10. With new insights about which features are relevant, iterate the process to step 2—or to step 1 if more test data is needed (using the models)—for subsequent analyses.

Vision : map the system – problem space

Lexin Zhou, Fernando Martínez-Plumed, José Hernández-Orallo Cèsar Ferri and Wout Schellaert “Reject Before You Run: Small Assessors Anticipate Big Language Models”
EBeM@IJCAI2022

- Identify dimensions in systems (capabilities) and problems (difficulties):
 - The assessor is a simple parametric model.
- Otherwise, use non-parametric assessors.

VISION:
Having **every deployed AI system backed by and accounted for with its capability profile and/or its assessor model**



Thank you!

JOSE H. ORALLO

<http://josephorallo.webs.upv.es/>
jorallo@upv.es

