# Common Task Framework For a Critical Evaluation of Scientific Machine Learning Algorithms

Philippe M. Wyder, Judah Goldfeder, Alexey Yermakov, Yue Zhao, Stefano Riva, Jan Williams, David Zoro, Amy Sara Rude, Matteo Tomasetto, Joe Germany, Joseph Bakarji, Georg Maierhofer, Miles Cranmer, J. Nathan Kutz

# No Common Standard

- The field is suffering from:

  - Weak Baselines

  - Reporting Bias

  - Inconsistent Evaluations
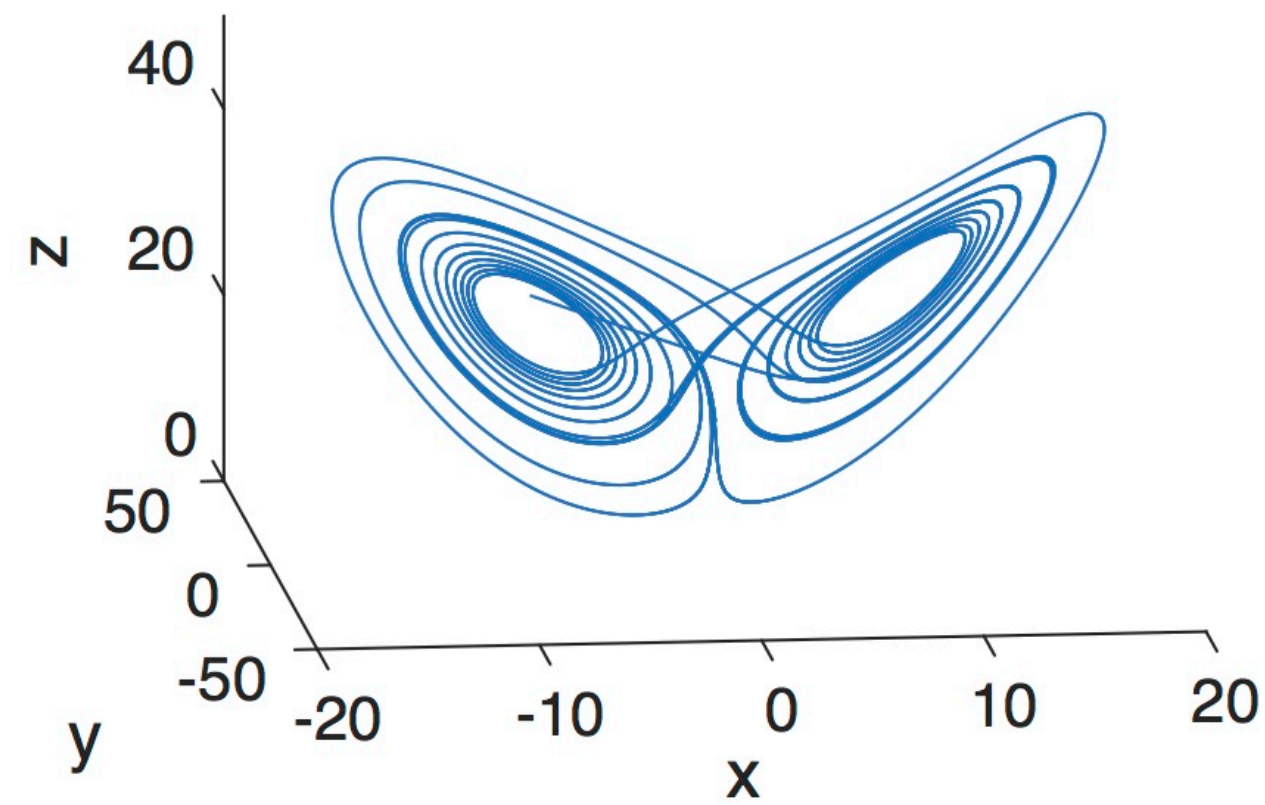
  - Lack of *hidden* test sets

# Common Task Framework for Science

- We provide:

  - Curated set of datasets

  - Task specific metrics

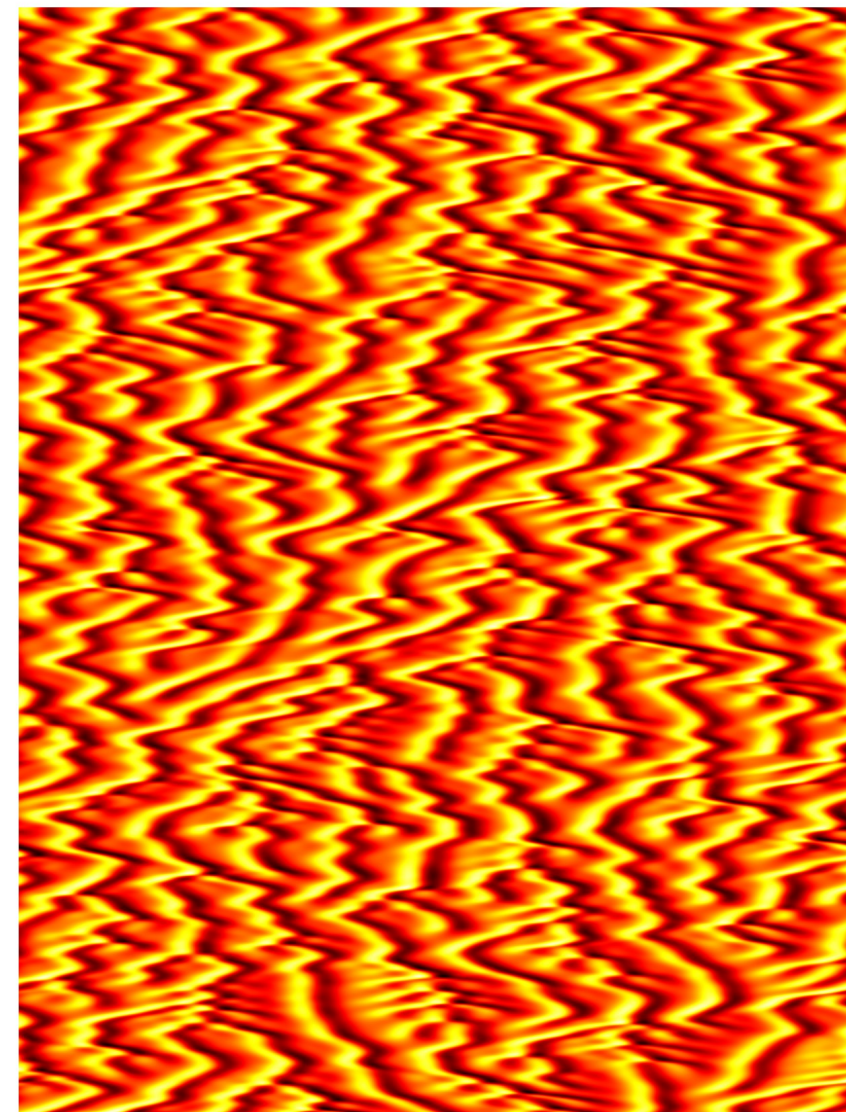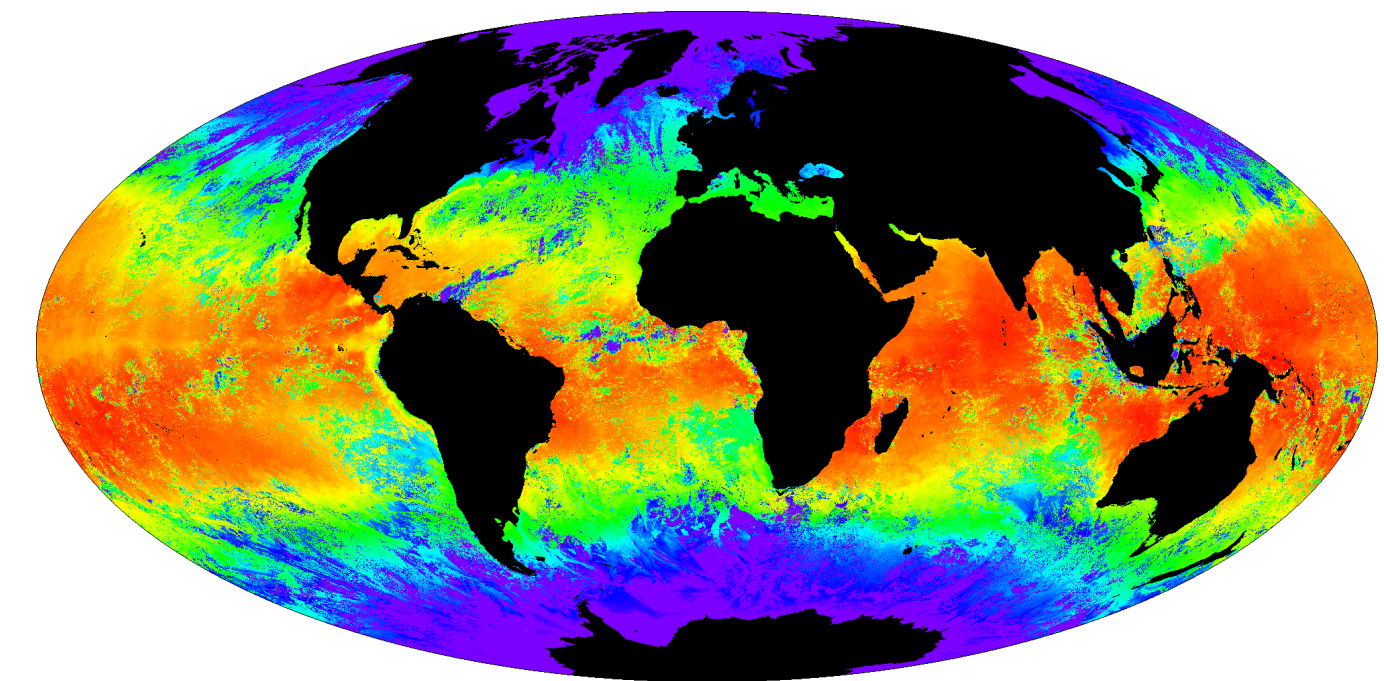  - Standardized Evaluation

  - Hidden test set

# Datasets

Lorenz

Kuramoto-Sivashinsky

Sea Surface Temperature



$$\mathbf{x}_j = \mathbf{x}(t_j) = \begin{bmatrix} x \\ y \\ z \end{bmatrix}$$

$$\mathbf{x}_j = \mathbf{x}(t_j) = \begin{bmatrix} u_1 \\ u_2 \\ \vdots \\ u_n \end{bmatrix}$$

More to come!

NEURAL INFORMATION
PROCESSING SYSTEMS

# Average Scores



Lorenz

| Model | Average Score |
|---|---|
| LSTM | 64.54 |
| DeepONet | 57.80 |
| Reservoir | 54.87 |
| ODE-LSTM | 41.67 |
| Spacetime | 38.79 |
| KAN | 35.32 |
| SINDy | 33.08 |
| Opt DMD | 20.62 |
| FNO | 20.51 |
| PyKoopman | 8.33 |
| NeuralODE | 3.38 |
| HigherOrder DMD | 0.94 |
| Baseline Average | -4.73 |
| Baseline Zeros | -39.00 |

Kuramoto-Sivashinsky

| Model | Average Score |
|---|---|
| Reservoir | 18.88 |
| LSTM | 15.61 |
| ODE-LSTM | 15.28 |
| DeepONet | 6.99 |
| KAN | 6.54 |
| Baseline Zeros | 0.00 |
| Baseline Average | -3.02 |
| SINDy | -5.23 |
| Opt DMD | -11.04 |
| FNO | -18.51 |
| PyKoopman | -20.11 |
| HigherOrder DMD | -24.98 |
| NeuralODE | -33.46 |
| Spacetime | -45.58 |

# Characterizing Model Performance
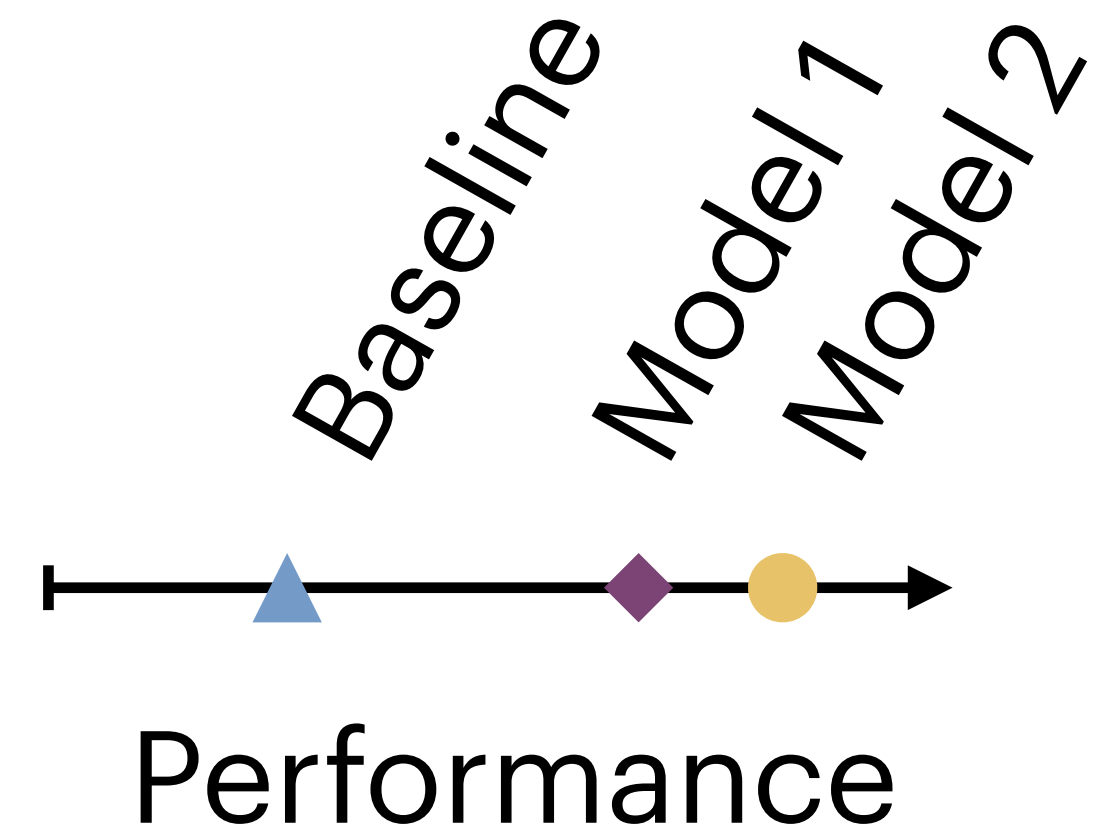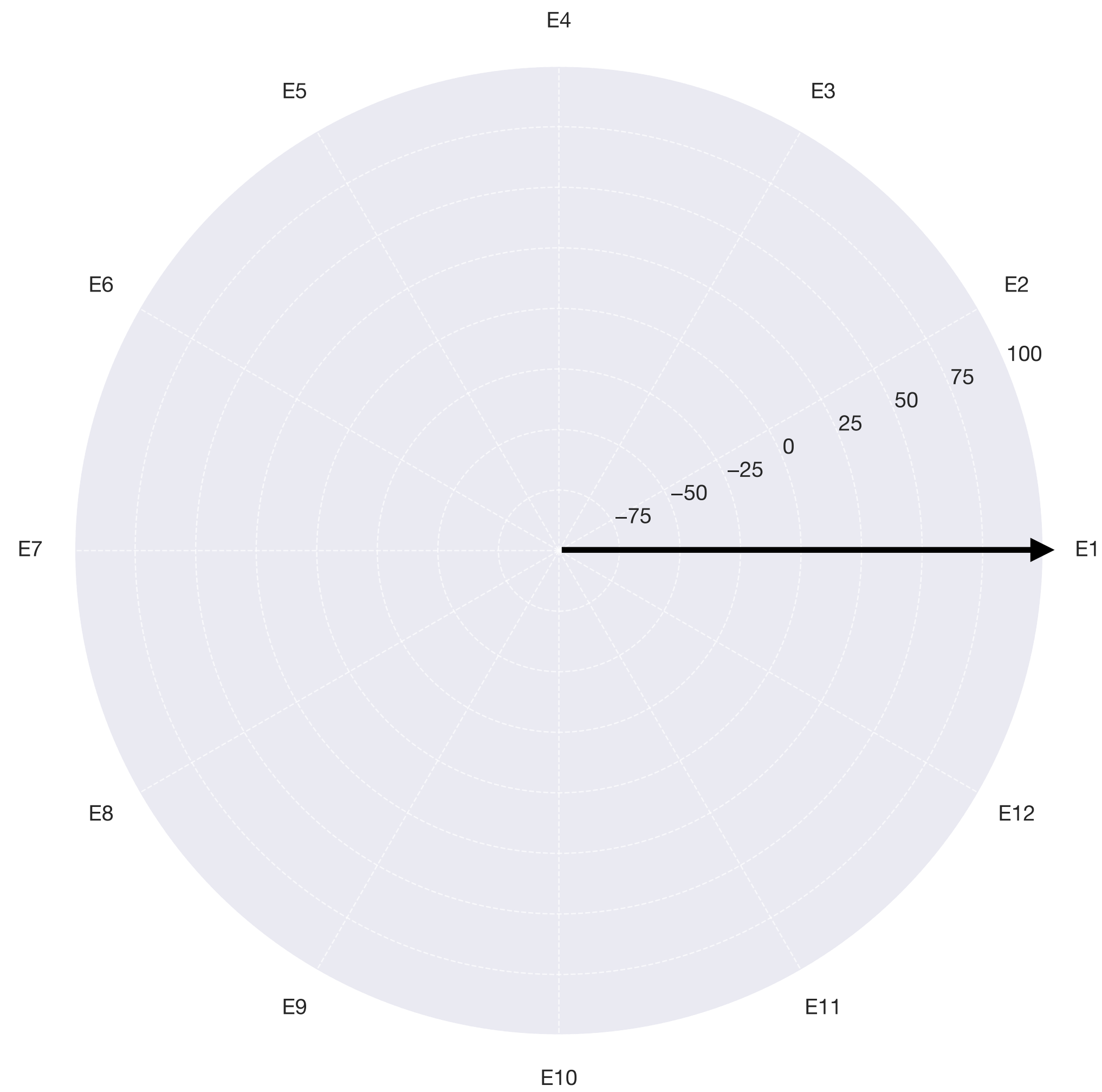
# Characterizing Model Performance

# Characterizing Model Performance



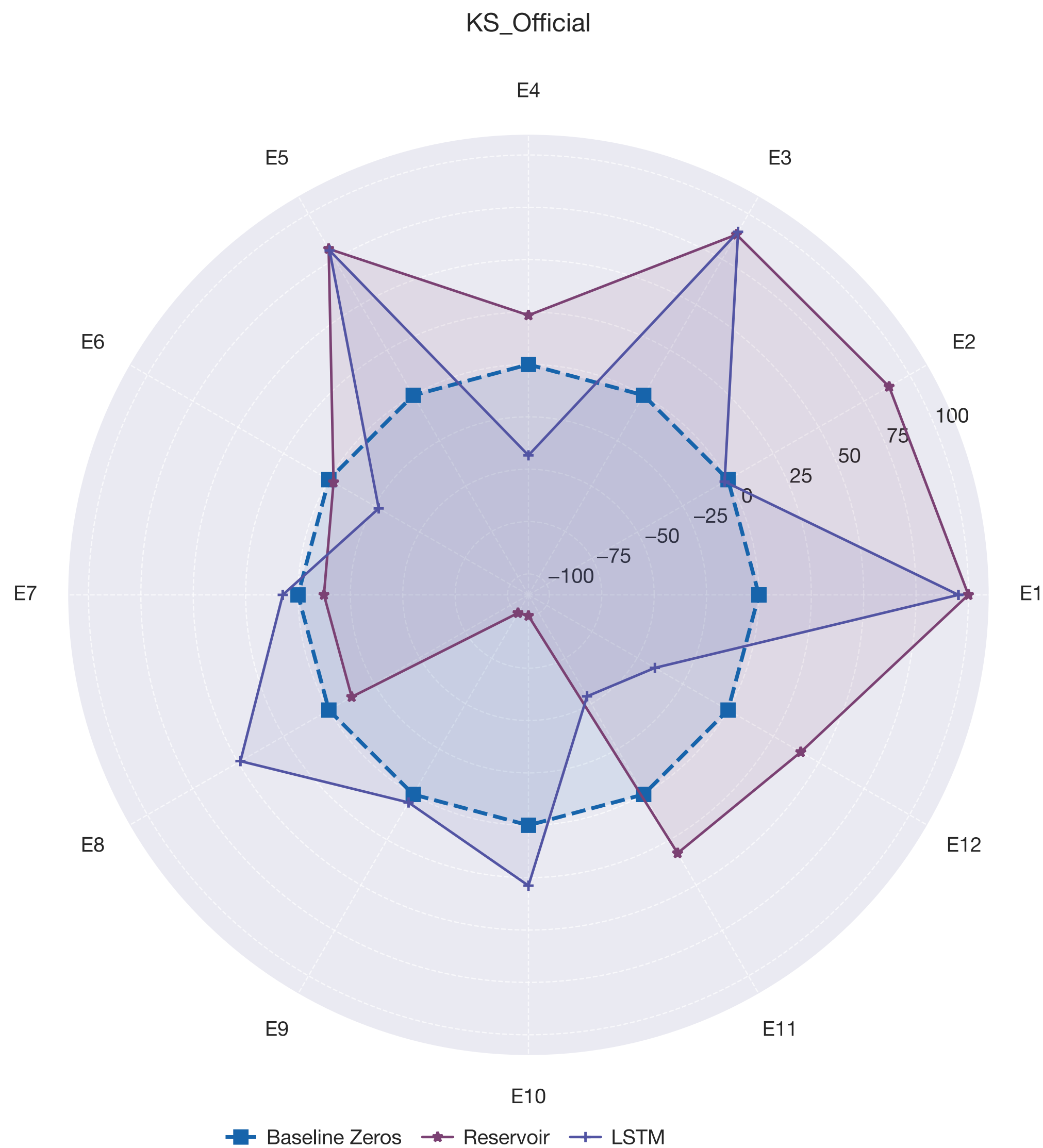Scoring Dimensions

| Score | Test | Task |
|-------|------|------|
| E1 | Forecasting | Short-time |
| E2 | Forecasting | Long-time |
| E3 | Noisy (medium) | Reconstruction (denoising) |
| E4 | Noisy (medium) | Forecast (long-time) |
| E5 | Noisy (high) | Reconstruction (denoising) |
| E6 | Noisy (high) | Forecast (long-time) |
| E7 | Limited Data (clean) | Forecast (short-time) |
| E8 | Limited Data (clean) | Forecast (long-time) |
| E9 | Limited Data (noisy) | Forecast (short-time) |
| E10 | Limited Data (noisy) | Forecast (long-time) |
| E11 | Parametric Generalization | Interpolation forecast |
| E12 | Parametric Generalization | Extrapolation forecast |

# Model Fingerprint

# CTF4Science on GitHub

- Modular code with diverse functionality

  - Data loading and model evaluation API

  - Easy hyper-parameter tuning

- Pre-configured models

  - 18 models implemented, tuned, and evaluated

- Visit https://github.com/CTF-for-Science/ctf4science

NEURAL INFORMATION
PROCESSING SYSTEMS

# CTF for Science on Kaggle

- Let the broader community score their models

- Truly withheld test-set for fair comparison

- Display submissions per team to prevent p-hacking

- Custom website to show spider plots and detailed metrics

NEURAL INFORMATION
PROCESSING SYSTEMS