# Benchmarking out-of-the-box forecasters of varying scales in biology

**Anthony Culos**
Department of Computer Science
Columbia University
New York, NY 10027
aec2244@columbia.edu

**Mohammed AlQuraishi**
Department of Computer Science
Department of Systems Biology
New York, NY 10027
ma4129@columbia.edu

## Abstract

Forecasting in biological systems presents different considerations and difficulties from traditional time-series settings, most notably the high-dimensionality associated with modern biological assays. A comprehensive analysis of the performance of modern forecasting methods on biological datasets is currently missing from the literature, in particular one that evaluates potential gains from larger model sizes and their increased complexity. Here, we assess 14 models spanning 4 complexity scales (Baseline, Statistical, Neural, and LLM-based) on 5 time-series datasets. We show that model scale and complexity does not uniformly improve performance across biological datasets, and that in some cases, highly complex models fail to outperform common baselines.

## 1 Introduction

Time series forecasting has historically been developed and applied by statisticians and econometricians. In recent years, the machine learning community has made strides in developing models that leverage increasing amounts of available longitudinal data. Forecasting problems native to biology however do not take advantage of pre-existing models, opting instead to develop bespoke predictive algorithms [9, 8]. Out-of-the-box forecasting models could reduce time spent creating highly specialized models allowing researchers to focus on understanding biological phenomena. Evaluation of general forecasting models in biological settings has, to our knowledge, not yet been rigorously studied: Of particular interest is the effect model size and complexity has on performance as larger and more elaborate models become available. As such, we evaluate 4 classic baselines, 3 statistical, 6 neural, and one LLM model across 5 domains: The 4 Electric Transformer Temperature (ETT) datasets, 3 selected Chaotic Attractors, 4 Microbiome datasets, 7 single cell Ribonucleic Acid (scRNA) sequencing datasets, and an epidemiological dataset. See table 1 in Appendix A for quantitative descriptions of each.

Non-biological datasets were chosen to give a baseline perspective on model performance, as the ETT and Chaotic Attractor datasets are better understood systems than their biological counterparts. The three biological datasets were selected for their diverse representation of biological scales and high dimensionality. While not a complete view of biology, they provide a sample of contexts in which forecasting models can be applied.

To give a robust assessment of model performance, we use multiple commonly used evaluation metrics, namely: Mean Absolute Error (MAE), Mean Squared Error (MSE), Mean Absolute Percentage Error (MAPE), and symmetric MAPE (sMAPE). While giving a qualitative description of results would help contextualize model performance, what constitutes *"good"* or *"bad"* in biology is highly varied and nuanced. For instance, a computational biologist may consider a small error relative to a benchmark model sufficiently good whereas a medical doctor may require that a model be accurate enough to

inform their diagnostic decisions. Consequently, we avoid making any qualitative claims for any specific dataset, instead opting only to discuss models in terms of the evaluation metrics and their performance across different datasets.

Since many datasets are non-standard we select a wide variety of models to generate as thorough an evaluation as possible. We include four classic baselines: the Naive forecast, a Random Walk with Drift (RWD), a Window Average, and the Historic Average. Three well established hyper-parameter optimized statistical models: the Auto Regressive Integrated Moving Average (ARIMA), Error-Trend-Seasonality (ETS), and Theta models. Six modern deep-learning based models: the Autoformer, Informer, N-Beats, N-HiTs, TSMixer, and TimeMixer models [20, 22, 14, 2, 3, 19]. Along with Time-LLM using the Gemma2-2B as the largest model we investigate [11, 18].

## 2    Methodology

### 2.1    Data collection & background

Chaotic Attractors have been proposed by Gilpin [10] as a benchmark for forecasting since they are well understood systems which are notoriously difficult to predict. There is also some evidence that chemical kinetics operate as a chaotic system making them a potential surrogate for cell-signalling in lieu of available data [17]. While Gilpin provides dozens of options, we arbitrarily selected 3 (Aizawa, Blasius, and Lorenz) to not remove the biological focus and better evaluate a smaller number of chaotic systems. ETT, the other non-biological dataset, measures the load and oil temperature of two electrical transformers at two stations from a single Chinese province [22].

scRNA seq data was originally collected to infer gene regulatory networks from single-cell transcriptomics [15]. Experimental data included is from human and mouse cell lines, specifically human Embryonic Stem Cells (hESC), human Hepatocyte-like cells (HEPs), mouse Embryonic Stem Cells (mESC), mouse Dendritic Cells (mDC), and mouse Hematopoietic Stem Cells (mHSC-E, mHSC-GM, mHSC-L) from three lineages Erythroid, Granulocyte-Monocyte, and Lymphoid like. While used in cell-fate prediction and gene-regulatory inference, these transcriptional measures do not explicitly include time. Instead a *"pseudotime"* is inferred using an ordering of cells by similarity of gene expression computed using Slingshot [16]. We filter the psuedotime measure to remove duplicates and convert it to date-time for use with general forecasters.

Microbiome data were collected from 4 individuals across two publications [1, 6], both studying the long term dynamics of microbial communities. Here each measurement represents the size of closely related microbial populations called Operational Taxonomic Units (OTUs) and the time-series represents a complex ecology of microorganisms.

Epidemiological data was aggregated across Chinese provinces following China's exit from strict zero-COVID policy for approximately three years by [21]. Measures focus on different respiratory and cardiological ailments aggregated across seven hospitals.

### 2.2    Model fitting

All data was split with a 70% train, 10% validation, and 20% test ratio. Classic baselines, statistical, and neural models were optimized and fit using the NIXTLA suite of tools [7, 13]. For Time-LLM we used its original implementation. MAE was used as the loss metric for fitting and optimizing deep learning based models except for Time-LLM which used MSE. Each dataset fit models to a short and a long horizon to also investigate sensitivity to prediction size. All models were fit with a maximum compute budget of 5 days on a single Nvidia A6000 GPU, except for Time-LLM which used 6 GPUs in order to fit the Gemma2-2B LLM in memory. Microbiome and scRNA data also frequently exceeded the practical maximum number of dimensions for application of TSMixer, TimeMixer, and Time-LLM models and as such data exceeding 524 dimensions were sub-sampled to include the highest variance 524 features.

## 3    Results

A total of 2,128 trained model evaluations across each dataset, model, prediction horizon, and evaluation metrics were recorded for both training and test data. To reduce complexity, we will focus

on test set MAE and sMAPE forecast evaluations for the longer time horizon task across all models and datasets (Fig 1). All evaluation metrics are calculated on predictions in the data's original scale prior to any pre-processing done by the model. If a model predicts a scaled forecast they are un-scaled prior to evaluation which may lead to larger reported metrics than previous publications.

Since MAE summarizes errors directly calculated from observations while sMAPE summarizes relative errors the two metrics can, at times, be discordant. For instance, many OTUs in the microbiome dataset exist at a level difficult to detect (ergo measured as zero) or are introduced to the microbial ecology rather disruptively, leading to long periods of dormancy followed by high levels of activity. This in turn can cause MAE values to accumulate errors from outliers more aggressively than sMAPE, as the latter is bounded by being a proportion, resulting in a more pessimistic MAE assessment of microbiome data for baseline and statistical models. A similar, but less pronounced, phenomenon is visible within the human scRNA sequencing data.

Improved model performance for Chaotic Attractors, ETT, and COVID was associated with an increase in model complexity while Microbiome and scRNA seq data failed to see such gains. This lack of difference in model performance for biological datasets could stem from similar data not being used during model development. Epidemiological data is an exception among biological datasets as its forecasting was very recently a matter of international safety and models perform well on this task. Alternatively the variables, exogenous or otherwise, needed to indicate the observed dynamics may not be available to the models.



Figure 1: Test set MAE and sMAPE model evaluation on the longest prediction horizon task for each dataset (Lower is better for both).

Inspecting a selection of models on the Lorenz chaotic attractor gives some insights into where these models may be struggling in the biological contexts [12]. Evident from Fig. 2, neural models tend to do well within the basin of each attractor but struggle in the transition between these two modes. Many biological systems of interest exhibit similar structural changes. The microbiome datasets have known non-stationarities, with donorA traveling abroad introducing new microbial communities and donorB experiencing salmonella. Such dynamics can cause rapid but temporary shifts to different modes for each time-series.

Figure 2: Test set prediction on Lorenz Attractor for a selection of models.

# 4   Discussion

COVID was the one biological dataset to exhibit performance characteristics similar to non-biological data. Its relatively large MAE is to be expected as respiratory diseases and related ailments measured were during a respiratory disease pandemic in a population of roughly 1.4 billion. More interesting is where model size and complexity failed to provide performance gains, as it could provide insights into future research directions. Forecasting in both microbiome and scRNA data fail to outperform the baseline models in most cases. These biological data offer a new test bed for future models to evaluate performance on what appears to be an out-of-distribution (ood) task compared to more standard datasets like ETT. We believe that the choice of benchmark datasets directly guide model development (ImageNet, CIFAR-100, and Enron Datasets for example). Inclusion of more biologically relevant time series could lead to more powerful general purpose forecasters in biology. Biological data and foundation models also present a novel intersection, both as an ood evaluation of foundation model performance and as diverse data sources. Complimentary to this, quantification of exogenous variables impact on performance in these biological forecasters could lead to novel discoveries of upstream dependencies of the systems' dynamics.

While foundation models offer an exciting research avenue, their potential for inadvertent data leakage given the large amount of training data raises some concern, especially in biological and medical 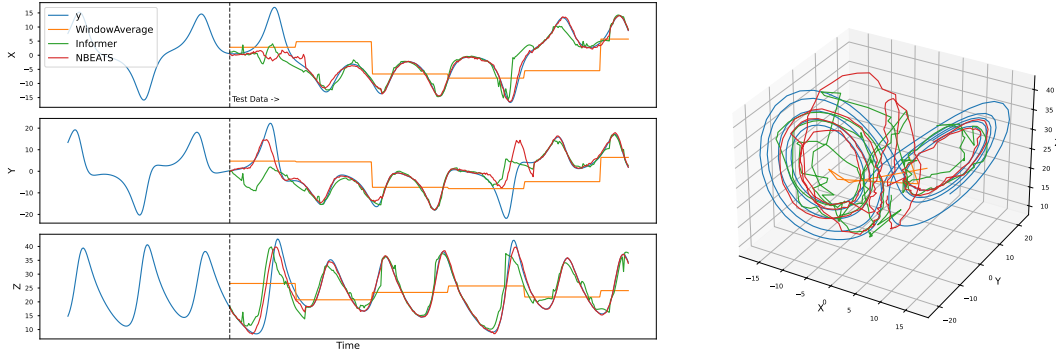applications. For instance, Time-LLM performance shows an improvement using the Gemma2-2B LLM. However, from these results alone we cannot conclusively attribute this gain to genuine model improvements from the LLama1-7B and GPT2 models as Gemma2-2B may have memorized ETT data (they were published prior to Gemma2-2B's publication). Extending this to biological and healthcare contexts, the potential for data leakage could falsely indicate improved understanding of biological principles or overestimate the performance of clinical predictors.

A noticeable omission from our biological datasets is cell signalling data, a prominent domain of study involving molecular interactions which facilitate cell-to-cell communication and gene regulation. This is due to the way in which cell-signalling data is collected, predomi-tly by flow or mass cytometry. The latter being a destructive measure of the cell making it difficult to track through time. Regardless of whether cell viability is maintained, measurements are often collected in a non-uniform manner as early dynamics tend to be more relevant to the specific phenomena being studied.

As biological assays continue to improve and decrease in cost, a proactive approach to high-dimensional model development would facilitate rapid application to biological domains with the potential for impactful results in both biology and healthcare. One technique of particular interest is LiveSEQ, it measures single cell transcription of genes while preserving their viability allowing for repeated measurements [5]. Models that can capture the whole genome (approximately 20 thousand genes) to facilitate a systems-level perspective of biology are of particular interest. Many biological systems presented here have classically been studied using extensions of basic dynamical systems, such as the Lotka-Voltaire predator prey model applied to microbial ecology. Application of Neural ODEs either as a standalone forecasting model or a component in a foundation time-series model may alleviate some of the issues native to biological data discussed here [4].

4

# References

[1] Caporaso, J. G., Lauber, C. L., Costello, E. K., Berg-Lyons, D., Gonzalez, A., Stombaugh, J., Knights, D., Gajer, P., Ravel, J., Fierer, N., et al. (2011). Moving pictures of the human microbiome. *Genome biology*, 12:1–8.

[2] Challu, C., Olivares, K. G., Oreshkin, B. N., Ramirez, F. G., Canseco, M. M., and Dubrawski, A. (2023). Nhits: Neural hierarchical interpolation for time series forecasting. In *Proceedings of the AAAI conference on artificial intelligence*, volume 37, pages 6989–6997.

[3] Chen, S.-A., Li, C.-L., Arik, S. O., Yoder, N. C., and Pfister, T. (2023). TSMixer: An all-MLP architecture for time series forecast-ing. *Transactions on Machine Learning Research*.

[4] Chen, T. Q., Rubanova, Y., Bettencourt, J., and Duvenaud, D. K. (2018). Neural ordinary differential equations. In *Neural Information Processing Systems*.

[5] Chen, W., Guillaume-Gentil, O., Rainer, P. Y., Gäbelein, C. G., Saelens, W., Gardeux, V., Klaeger, A., Dainese, R., Zachara, M., Zambelli, T., Vorholt, J. A., and Deplancke, B. (2022). Live-seq enables temporal transcriptomic recording of single cells. *Nature*, 608(7924):733–740.

[6] David, L. A., Materna, A. C., Friedman, J., Campos-Baptista, M. I., Blackburn, M. C., Perrotta, A., Erdman, S. E., and Alm, E. J. (2014). Host lifestyle affects human microbiota on daily timescales. *Genome biology*, 15:1–15.

[7] Federico Garza, Max Mergenthaler Canseco, C. C. K. G. O. (2022). StatsForecast: Lightning fast forecasting with statistical and econometric models. PyCon Salt Lake City, Utah, US 2022.

[8] Gabor, A., Tognetti, M., Driessen, A., Tanevski, J., Guo, B., Cao, W., Shen, H., Yu, T., Chung, V., in Breast Cancer DREAM Consortium members, S. C. S., et al. (2021). Cell-to-cell and type-to-type heterogeneity of signaling networks: insights from the crowd. *Molecular systems biology*, 17(10):e10402.

[9] Gerber, G. K. (2014). The dynamic microbiome. *FEBS letters*, 588(22):4131–4139.

[10] Gilpin, W. (2021). Chaos as an interpretable benchmark for forecasting and data-driven modelling. In *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 2)*.

[11] Jin, M., Wang, S., Ma, L., Chu, Z., Zhang, J. Y., Shi, X., Chen, P.-Y., Liang, Y., Li, Y.-F., Pan, S., and Wen, Q. (2024). Time-LLM: Time series forecasting by reprogramming large language models. In *The Twelfth International Conference on Learning Representations*.

[12] Lorenz, E. N. (1963). Deterministic nonperiodic flow. *Journal of atmospheric sciences*, 20(2):130–141.

[13] Olivares, K. G., Challú, C., Garza, F., Canseco, M. M., and Dubrawski, A. (2022). NeuralForecast: User friendly state-of-the-art neural forecasting models. PyCon Salt Lake City, Utah, US 2022.

[14] Oreshkin, B. N., Carpov, D., Chapados, N., and Bengio, Y. (2020). N-beats: Neural basis expansion analysis for interpretable time series forecasting. In *International Conference on Learning Representations*.

[15] Pratapa, A., Jalihal, A. P., Law, J. N., Bharadwaj, A., and Murali, T. M. (2020). Benchmarking algorithms for gene regulatory network inference from single-cell transcriptomic data. *Nature Methods*, 17(2):147–154.

[16] Street, K., Risso, D., Fletcher, R. B., Das, D., Ngai, J., Yosef, N., Purdom, E., and Dudoit, S. (2018). Slingshot: cell lineage and pseudotime inference for single-cell transcriptomics. *BMC Genomics*, 19(1):477.

[17] Susits, M. and Tóth, J. (2024). Rigorously proven chaos in chemical kinetics.

[18] Team, G., Riviere, M., Pathak, S., Sessa, P. G., Hardin, C., Bhupatiraju, S., Hussenot, L., Mesnard, T., Shahriari, B., Ramé, A., et al. (2024). Gemma 2: Improving open language models at a practical size. *arXiv preprint arXiv:2408.00118*.

[19] Wang, S., Wu, H., Shi, X., Hu, T., Luo, H., Ma, L., Zhang, J. Y., and ZHOU, J. (2024). Timemixer: Decomposable multiscale mixing for time series forecasting. In *The Twelfth International Conference on Learning Representations*.

[20] Wu, H., Xu, J., Wang, J., and Long, M. (2021). Autoformer: Decomposition transformers with auto-correlation for long-term series forecasting. *Advances in neural information processing systems*, 34:22419–22430.

[21] Yu, X.-s., Tan, S., Tang, W., Zhao, F.-f., Ji, J., Lin, J., He, H.-j., Gu, Y., Liang, J.-J., Wang, M., Chen, Y., Yang, J., Xie, L., Wang, Q., Liu, M., He, Y., Chen, L., Wang, Y. X., Wu, Z., Zhao, G., Liu, Y., Wang, Y., Hao, D., Cen, J., Yao, S.-Q., Zhang, D., Liu, L., Lye, D. C., Hao, Z., Wong, T. Y., and Cen, L.-P. (2024). Multi-dimensional epidemiology and informatics data on covid-19 wave at the end of zero covid policy in china. *Frontiers in Public Health*, 12.

[22] Zhou, H., Zhang, S., Peng, J., Zhang, S., Li, J., Xiong, H., and Zhang, W. (2021). Informer: Beyond efficient transformer for long sequence time-series forecasting. In *Proceedings of the AAAI conference on artificial intelligence*, volume 35, pages 11106–11115.

# A Appendix

## A.1 Data availability

All data and scripts for reproduction of results can be found at https://figshare.com/s/34a486f130bb8c4aee2b.

## A.2 Supplemental tables

| Dataset | Dimension | Total Trajectory Length | Long Horizon | Short Horizon | Frequency |
|---|---|---|---|---|---|
| All Dynamical Systems | 3 | 2000 | 64 | 32 | Second |
| ETTh1/ETTh2 | 7 | 14400 | 192 | 96 | Hour |
| ETTm1/ETTm2 | 7 | 57600 | 192 | 96 | 15 Min |
| COVID | 36 | 955 | 30 | 14 | Day |
| Microbiome - donorA | 524 (1570) | 365 | 8 | 4 | Day |
| Microbiome - donorB | 524 (1525) | 253 | 8 | 4 | Day |
| Microbiome - female | 524 (552) | 185 | 8 | 4 | Day |
| Microbiome - male | 524 (1254) | 443 | 8 | 4 | Day |
| scRNA - hESC | 524 (17745) | 758 | 16 | 8 | Second |
| scRNA - hHep | 524 (11515) | 425 | 16 | 8 | Second |
| scRNA - mDC | 524 (7371) | 383 | 16 | 8 | Second |
| scRNA - mESC | 524 (18385) | 421 | 16 | 8 | Second |
| scRNA - mHSC-L | 524 (4762) | 847 | 16 | 8 | Second |
| scRNA -mHSC-E | 524 (4762) | 1071 | 16 | 8 | Second |
| scRNA -mHSC-GM | 524 (4762) | 889 | 16 | 8 | Second |

Table 1: Summary of datasets used for this work, note that when subsampling was done actual dimensionality is shown in parenthesis.

| | | Naive | | RWD | | Average | | Window Average | | ARIMA | | ETS | | THETA | | Autoformer | | Informer | | N-BEATS | | N-HITS | | TSMixer | | Time Mixer | | Time-LLM | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Dataset (H) | MAE | MSE | MAE | MSE | MAE | MSE | MAE | MSE | MAE | MSE | MAE | MSE | MAE | MSE | MAE | MSE | MAE | MSE | MAE | MSE | MAE | MSE | MAE | MSE | MAE | MSE | MAE | MSE |
| dysf | Aizawa (32) | 0.415 | 0.394 | 0.618 | 1.040 | 0.574 | 0.556 | 0.574 | 0.556 | 0.373 | 0.332 | 0.365 | 0.436 | 0.449 | 0.546 | 0.137 | 0.036 | 0.037 | 0.002 | 0.007 | 0.002 | 0.008 | 0.000 | 0.036 | 0.003 | 0.020 | 0.001 | 0.330 | 0.239 |
| dysf | Aizawa (64) | 0.555 | 0.569 | 0.656 | 0.842 | 0.580 | 0.594 | 0.580 | 0.594 | 0.519 | 0.508 | 0.695 | 1.117 | 0.597 | 0.681 | 0.269 | 0.185 | 0.041 | 0.004 | 0.015 | 0.002 | 0.016 | 0.002 | 0.171 | 0.071 | 0.040 | 0.004 | - | - |
| dysf | Blasius (32) | 1.542 | 6.157 | 2.243 | 13.880 | 2.264 | 10.500 | 2.264 | 10.500 | 1.419 | 5.461 | 1.242 | 5.348 | 1.460 | 5.737 | 1.171 | 2.524 | 0.351 | 0.412 | 0.058 | 0.024 | 0.066 | 0.098 | 0.268 | 0.208 | 0.114 | 0.039 | 0.553 | 1.130 |
| dysf | Blasius (64) | 2.250 | 11.346 | 3.529 | 29.865 | 1.909 | 7.498 | 1.909 | 7.498 | 2.152 | 10.698 | 2.509 | 16.254 | 2.227 | 11.316 | 1.627 | 5.752 | 0.310 | 0.282 | 0.042 | 0.006 | 0.076 | 0.045 | 0.155 | 0.056 | 0.097 | 0.026 | - | - |
| dysf | Lorenz (32) | 7.100 | 88.865 | 10.454 | 201.171 | 7.342 | 86.320 | 7.342 | 86.320 | 6.189 | 78.601 | 7.767 | 128.536 | 7.230 | 97.280 | 4.393 | 41.826 | 1.652 | 8.374 | 0.342 | 0.764 | 0.395 | 0.935 | 3.674 | 44.981 | 1.069 | 3.727 | 4.795 | 42.023 |
| dysf | Lorenz (64) | 6.986 | 85.978 | 10.222 | 175.638 | 7.074 | 72.860 | 7.074 | 72.860 | 6.839 | 90.557 | 9.251 | 155.344 | 7.064 | 90.819 | 4.372 | 41.802 | 1.533 | 6.653 | 0.569 | 1.392 | 0.341 | 0.585 | 1.868 | 9.101 | 0.594 | 1.121 | - | - |
| ETT | ETTh1 (96) | 0.583 | 0.725 | 0.738 | 1.082 | 0.541 | 0.567 | 0.541 | 0.567 | 0.508 | 0.518 | 0.572 | 0.708 | 0.593 | 0.742 | 0.471 | 0.466 | 0.446 | 0.442 | 0.399 | 0.359 | 0.402 | 0.365 | 0.411 | 0.390 | 0.412 | 0.390 | 1.311 | 5.092 |
| ETT | ETTh1 (192) | 0.635 | 0.845 | 0.809 | 1.370 | 0.568 | 0.620 | 0.568 | 0.620 | 0.561 | 0.625 | 0.626 | 0.826 | 0.652 | 0.910 | 0.509 | 0.565 | 0.477 | 0.498 | 0.440 | 0.423 | 0.434 | 0.412 | 0.446 | 0.453 | 0.456 | 0.467 | - | - |
| ETT | ETTh2 (96) | 0.391 | 0.497 | 0.538 | 1.076 | 0.383 | 0.402 | 0.383 | 0.402 | 0.355 | 0.346 | 0.382 | 0.486 | 0.429 | 0.680 | 0.368 | 0.358 | 0.358 | 0.378 | 0.294 | 0.301 | 0.304 | 0.323 | 0.267 | 0.210 | 0.300 | 0.242 | 2.509 | 20.504 |
| ETT | ETTh2 (192) | 0.449 | 0.617 | 0.585 | 1.042 | 0.442 | 0.541 | 0.442 | 0.541 | 0.406 | 0.449 | 0.444 | 0.609 | 0.482 | 0.686 | 0.438 | 0.524 | 0.182 | 0.094 | 0.358 | 0.524 | 0.359 | 0.398 | 0.365 | 0.396 | 5.937e7 | 1.281e16 | - | - |
| ETT | ETTm1 (96) | 0.590 | 0.829 | 0.768 | 1.397 | 0.515 | 0.541 | 0.515 | 0.541 | 0.576 | 0.811 | 0.628 | 0.925 | 0.587 | 0.848 | 0.397 | 0.358 | 0.302 | 0.233 | 0.261 | 0.193 | 0.256 | 0.186 | 0.289 | 0.228 | 0.293 | 0.233 | 1.161 | 4.170 |
| ETT | ETTm1 (192) | 0.653 | 0.942 | 0.791 | 1.357 | 0.541 | 0.585 | 0.541 | 0.585 | 0.636 | 0.921 | 0.693 | 1.048 | 0.641 | 0.921 | 0.452 | 0.438 | 0.387 | 0.363 | 0.354 | 0.363 | 0.360 | 0.327 | 0.334 | 0.273 | 0.282 | 0.197 | - | - |
| ETT | ETTm2 (96) | 0.306 | 0.252 | 0.391 | 0.436 | 0.314 | 0.268 | 0.314 | 0.268 | 0.297 | 0.239 | 0.305 | 0.254 | 0.316 | 0.283 | 0.299 | 0.250 | 0.193 | 0.117 | 0.238 | 0.183 | 0.241 | 0.190 | 0.243 | 0.182 | 0.233 | 0.169 | 1.757 | 9.324 |
| ETT | ETTm2 (192) | 0.352 | 0.361 | 0.434 | 0.578 | 0.360 | 0.373 | 0.360 | 0.373 | 0.344 | 0.354 | 0.352 | 0.364 | 0.365 | 0.400 | 0.336 | 0.350 | 0.300 | 0.291 | 0.273 | 0.266 | 0.275 | 0.263 | 0.264 | 0.204 | 0.289 | 0.264 | - | - |
| micro | Male (4) | 19.456 | 1.739e4 | 32.514 | 4.479e4 | 19.947 | 1.601e4 | 19.947 | 1.601e4 | 14.973 | 1.141e4 | 18.931 | 1.459e4 | 20.249 | 1.724e4 | 19.469 | 1.564e4 | 18.713 | 1.640e4 | 18.127 | 1.610e4 | 17.934 | 1.582e4 | 12.530 | 1.255e4 | 13.980 | 1.614e4 | 34.233 | 2.937e4 |
| micro | Male (8) | 21.502 | 2.010e4 | 33.508 | 4.623e4 | 20.917 | 1.600e4 | 20.917 | 1.600e4 | 17.467 | 1.397e4 | 21.482 | 1.742e4 | 23.588 | 2.077e4 | 18.710 | 1.605e4 | 18.093 | 1.593e4 | 3.294e3 | 7.705e8 | 18.313 | 1.621e4 | 19.104 | 1.401e4 | 18.714 | 1.526e4 | - | - |
| micro | Female (4) | 12.936 | 1.231e4 | 20.063 | 3.266e4 | 13.660 | 1.275e4 | 13.660 | 1.275e4 | 10.266 | 8.786e3 | 15.236 | 1.768e4 | 12.810 | 1.197e4 | 31.446 | 6.540e4 | 12.560 | 1.214e4 | 1.702e7 | 4.150e16 | 12.495 | 1.202e4 | 8.753 | 6.979e3 | 9.456 | 9.148e3 | 40.260 | 3.549e4 |
| micro | Female (8) | 15.146 | 1.610e4 | 24.733 | 4.472e4 | 14.351 | 1.235e4 | 14.351 | 1.235e4 | 11.694 | 9.506e3 | 15.458 | 1.550e4 | 16.349 | 1.705e4 | 13.269 | 1.108e4 | 13.425 | 1.230e4 | 1.775e4 | 2.816e13 | 15.124 | 1.502e4 | 7.234 | 8.128e3 | 7.891 | 4.184e3 | - | - |
| micro | Donor A (4) | 19.112 | 2.100e4 | 30.489 | 5.132e4 | 18.893 | 1.917e4 | 18.893 | 1.917e4 | 14.508 | 1.347e4 | 15.909 | 1.505e4 | 19.741 | 2.683e4 | 18.501 | 1.965e4 | 34.392 | 7.241e4 | 17.535 | 1.924e4 | 20.551 | 2.289e4 | 10.177 | 6.388e3 | 12.911 | 1.190e4 | 34.045 | 3.575e4 |
| micro | Donor A (8) | 20.523 | 2.230e4 | 29.594 | 4.349e4 | 18.893 | 2.094e4 | 18.893 | 2.094e4 | 17.000 | 1.710e4 | 18.507 | 1.840e4 | 21.635 | 2.648e4 | 19.148 | 2.018e4 | 18.619 | 2.114e4 | 17.927 | 2.005e4 | 1.583e3 | 4.924e8 | 9.738 | 5.515e3 | 6.274 | 2.027e3 | - | - |
| micro | Donor B (4) | 16.890 | 1.958e4 | 26.894 | 4.792e4 | 16.744 | 1.720e4 | 16.744 | 1.720e4 | 12.393 | 1.145e4 | 12.881 | 1.162e4 | 17.385 | 2.688e4 | 17.797 | 1.842e4 | 15.151 | 1.670e4 | 14.805 | 1.550e4 | 15.173 | 1.710e4 | 14.815 | 1.519e3 | 11.311 | 1.019e4 | 35.198 | 4.642e4 |
| micro | Donor B (8) | 19.214 | 2.395e4 | 28.773 | 5.462e4 | 17.377 | 1.744e4 | 17.377 | 1.744e4 | 14.742 | 1.427e4 | 15.727 | 1.605e4 | 20.461 | 3.313e4 | 16.959 | 1.841e4 | 15.795 | 1.824e4 | 16.616 | 2.035e4 | 16.333 | 1.943e4 | 16.131 | 1.805e4 | 11.236 | 1.060e4 | - | - |
| epi | Covid (14) | 322.247 | 9.991e5 | 453.470 | 2.112e6 | 351.586 | 1.591e6 | 351.586 | 1.591e6 | 363.481 | 1.180e6 | 1.607e4 | 1.100e10 | 364.407 | 9.736e5 | 465.649 | 1.545e7 | 475.904 | 2.255e7 | 428.269 | 1.385e7 | 444.656 | 1.385e7 | 480.304 | 1.584e7 | 269.544 | 3.884e6 | 582.367 | 1.977e7 |
| epi | Covid (30) | 350.291 | 1.174e6 | 496.837 | 3.208e6 | 421.765 | 1.954e6 | 421.765 | 1.954e6 | 492.947 | 2.284e6 | 1.496e3 | 4.091e7 | 364.569 | 1.229e6 | 711.569 | 2.561e7 | 407.477 | 1.986e7 | 17.865 | 3.138e4 | 522.513 | 1.909e7 | 717.740 | 2.574e7 | 312.109 | 5.388e6 | - | - |
| scRNA | hESC (8) | 2.589 | 12.593 | 3.854 | 27.165 | 2.107 | 7.164 | 2.107 | 7.164 | 2.066 | 6.931 | 2.275 | 7.765 | 2.457 | 10.375 | 2.165 | 7.738 | 2.132 | 7.957 | 2.117 | 7.913 | 2.127 | 7.983 | 2.125 | 7.617 | 2.013 | 7.283 | 2.125 | 6.905 |
| scRNA | hESC (16) | 2.655 | 13.047 | 3.804 | 26.137 | 2.117 | 7.088 | 2.117 | 7.088 | 2.104 | 7.005 | 2.167 | 7.063 | 2.285 | 8.605 | 2.140 | 7.868 | 2.134 | 8.007 | 2.076 | 7.621 | 2.109 | 7.715 | 1.996 | 6.998 | 2.084 | 7.438 | - | - |
| scRNA | hHep (8) | 2.462 | 11.769 | 3.651 | 24.653 | 2.065 | 6.903 | 2.065 | 6.903 | 2.032 | 6.713 | 2.108 | 6.725 | 2.357 | 9.489 | 1.944 | 6.883 | 1.927 | 6.937 | 1.933 | 6.927 | 1.935 | 6.964 | 1.925 | 6.661 | 1.461 | 4.412 | 2.014 | 6.142 |
| scRNA | hHep (16) | 2.668 | 13.455 | 3.807 | 25.941 | 2.138 | 7.197 | 2.138 | 7.197 | 2.151 | 7.398 | 2.159 | 7.189 | 2.310 | 8.841 | 1.998 | 7.416 | 1.979 | 7.455 | 1.973 | 7.394 | 1.980 | 7.280 | 1.580 | 4.832 | 1.855 | 6.141 | - | - |
| scRNA | mDC (8) | 1.691 | 5.189 | 2.679 | 12.603 | 1.306 | 2.759 | 1.306 | 2.759 | 1.306 | 2.788 | 1.301 | 3.049 | 1.698 | 5.083 | 1.296 | 2.785 | 1.289 | 2.873 | 1.237 | 2.703 | 1.304 | 2.929 | 0.680 | 1.049 | 0.878 | 1.570 | 1.149 | 2.172 |
| scRNA | mDC (16) | 1.767 | 5.654 | 2.392 | 9.983 | 1.294 | 2.690 | 1.294 | 2.690 | 1.324 | 2.831 | 1.281 | 2.853 | 1.429 | 3.383 | 1.282 | 2.755 | 1.278 | 2.881 | 1.290 | 2.881 | 1.297 | 3.027 | 1.266 | 2.686 | 0.932 | 1.737 | - | - |
| scRNA | mESC (8) | 0.831 | 1.370 | 1.205 | 2.898 | 0.669 | 0.818 | 0.669 | 0.818 | 0.659 | 0.809 | 0.691 | 0.848 | 0.790 | 1.194 | 0.669 | 0.836 | 0.664 | 0.843 | 0.666 | 0.843 | 0.665 | 0.849 | 0.591 | 0.700 | 0.651 | 0.825 | 0.662 | 0.789 |
| scRNA | mESC (16) | 0.841 | 1.397 | 1.161 | 2.692 | 0.675 | 0.819 | 0.675 | 0.819 | 0.675 | 0.842 | 0.669 | 0.812 | 0.733 | 1.004 | 0.656 | 0.814 | 0.670 | 0.874 | 0.662 | 0.825 | 0.661 | 0.817 | 0.657 | 0.825 | 0.609 | 0.740 | - | - |
| scRNA | mHSC-E (8) | 3.413 | 19.449 | 5.164 | 42.423 | 2.915 | 11.698 | 2.915 | 11.698 | 2.811 | 12.368 | 2.868 | 0.640 | 3.113 | 14.702 | 2.865 | 1.895 | 2.856 | 11.936 | 2.861 | 12.093 | 2.854 | 11.809 | 2.852 | 12.026 | 2.857 | 11.986 | 2.904 | 10.555 |
| scRNA | mHSC-E (16) | 3.471 | 19.976 | 5.108 | 41.168 | 3.027 | 12.069 | 3.027 | 12.069 | 2.919 | 12.069 | 2.938 | 1.219 | 3.127 | 14.398 | 2.944 | 12.271 | 2.908 | 11.779 | 2.904 | 12.147 | 2.933 | 12.340 | 2.727 | 10.861 | 2.744 | 11.179 | - | - |
| scRNA | mHSC-GM (8) | 3.484 | 19.803 | 5.216 | 42.801 | 2.897 | 11.372 | 2.897 | 11.372 | 2.801 | 11.656 | 2.841 | 0.220 | 3.224 | 15.802 | 2.838 | 11.994 | 2.823 | 11.752 | 2.819 | 11.909 | 2.826 | 11.790 | 2.633 | 10.568 | 2.807 | 11.856 | 2.813 | 10.132 |
| scRNA | mHSC-GM (16) | 3.493 | 20.029 | 5.015 | 39.717 | 2.903 | 11.021 | 2.903 | 11.021 | 2.852 | 12.466 | 2.850 | 0.278 | 3.097 | 13.848 | 2.841 | 11.794 | 2.896 | 11.571 | 2.821 | 11.560 | 2.832 | 11.145 | 2.621 | 10.521 | 2.779 | 11.757 | - | - |
| scRNA | mHSC-L (8) | 3.574 | 20.970 | 5.371 | 45.418 | 2.970 | 11.768 | 2.970 | 11.768 | 2.874 | 12.184 | 2.913 | 0.580 | 3.352 | 17.093 | 2.906 | 12.375 | 2.896 | 12.213 | 2.896 | 12.590 | 2.890 | 12.315 | 2.872 | 12.468 | 2.881 | 12.349 | 2.872 | 10.471 |
| scRNA | mHSC-L (16) | 3.566 | 20.885 | 5.113 | 41.150 | 2.943 | 11.193 | 2.943 | 11.193 | 2.910 | 12.955 | 2.913 | 0.558 | 3.159 | 14.240 | 2.876 | 11.914 | 2.861 | 12.397 | 2.869 | 12.304 | 2.862 | 11.987 | 2.754 | 11.619 | 2.731 | 11.638 | - | - |

Table 2: Train set MAE and MSE evaluations

Table 3: Test set MAE and MSE evaluations

| Group | Dataset (H) | Naive | | RWD | | Average | | Window Average | | ARIMA | | THETA | | ETS | | Autoformer | | Informer | | N-BEATS | | N-HiTS | | TSMixer | | Time Mixer | | Time-LLM | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | MAE | MSE | MAE | MSE | MAE | MSE | MAE | MSE | MAE | MSE | MAE | MSE | MAE | MSE | MAE | MSE | MAE | MSE | MAE | MSE | MAE | MSE | MAE | MSE | MAE | MSE | MAE | MSE |
| dysr | Aizawa (32) | 0.438 | 0.419 | 0.686 | 1.213 | 0.605 | 0.603 | 0.605 | 0.603 | 0.393 | 0.350 | 0.490 | 0.611 | 0.379 | 0.481 | 0.128 | 0.028 | 0.040 | 0.003 | 0.009 | 0.000 | 0.009 | 0.000 | 0.039 | 0.003 | 0.026 | 0.002 | - | - |
| | Aizawa (64) | 0.550 | 0.544 | 0.654 | 0.827 | 0.554 | 0.553 | 0.554 | 0.553 | 0.511 | 0.481 | 0.589 | 0.639 | 0.717 | 1.168 | 0.284 | 0.215 | 0.044 | 0.005 | 0.017 | 0.001 | 0.017 | 0.001 | 0.164 | 0.063 | 0.037 | 0.002 | 0.345 | 0.272 |
| | Blasius (32) | 1.609 | 6.387 | 2.347 | 14.284 | 2.492 | 11.705 | 2.492 | 11.705 | 1.478 | 5.638 | 1.525 | 5.945 | 1.299 | 5.253 | 1.049 | 2.053 | 0.331 | 0.282 | 0.037 | 0.005 | 0.127 | 0.083 | 0.245 | 0.164 | 0.126 | 0.036 | - | - |
| | Blasius (64) | 2.291 | 11.748 | 3.689 | 31.630 | 2.016 | 7.897 | 2.016 | 7.897 | 2.192 | 11.109 | 2.264 | 11.707 | 2.531 | 16.173 | 1.778 | 6.790 | 0.336 | 0.324 | 0.032 | 0.004 | 0.054 | 0.009 | 0.288 | 0.323 | 0.076 | 0.012 | 0.504 | 0.903 |
| | Lorenz (32) | 8.189 | 109.971 | 12.100 | 264.859 | 8.537 | 104.481 | 8.537 | 104.481 | 6.852 | 90.750 | 8.302 | 120.661 | 8.774 | 159.925 | 5.429 | 58.674 | 1.905 | 9.546 | 0.328 | 0.314 | 0.368 | 0.390 | 4.085 | 50.057 | 1.354 | 5.416 | 5.456 | 53.905 |
| | Lorenz (64) | 8.457 | 125.261 | 11.298 | 258.678 | 7.584 | 82.592 | 7.584 | 82.592 | 7.777 | 114.037 | 8.557 | 133.709 | 11.161 | 249.009 | 4.580 | 43.391 | 2.524 | 22.978 | 1.353 | 10.914 | 0.882 | 7.265 | 2.679 | 22.891 | 1.019 | 3.359 | - | - |
| ETT | ETTh1 (96) | 0.607 | 0.963 | 0.683 | 1.109 | 0.556 | 0.693 | 0.556 | 0.693 | 0.541 | 0.697 | 0.603 | 0.955 | 0.592 | 0.938 | 0.420 | 0.410 | 0.413 | 0.442 | 0.369 | 0.337 | 0.377 | 0.350 | 0.362 | 0.338 | 0.375 | 0.348 | - | - |
| | ETTh1 (192) | 0.602 | 0.895 | 0.711 | 1.182 | 0.571 | 0.701 | 0.571 | 0.701 | 0.557 | 0.715 | 0.608 | 0.925 | 0.592 | 0.883 | 0.421 | 0.408 | 0.430 | 0.423 | 0.403 | 0.376 | 0.397 | 0.368 | 0.394 | 0.362 | 0.410 | 0.383 | 1.148 | 5.388 |
| | ETTh2 (96) | 0.375 | 0.325 | 0.532 | 0.681 | 0.353 | 0.281 | 0.353 | 0.281 | 0.342 | 0.258 | 0.414 | 0.394 | 0.368 | 0.304 | 0.375 | 0.304 | 0.352 | 0.292 | 0.309 | 0.234 | 0.311 | 0.241 | 0.320 | 0.247 | 0.314 | 0.236 | - | - |
| | ETTh2 (192) | 0.433 | 0.433 | 0.648 | 0.945 | 0.404 | 0.354 | 0.404 | 0.354 | 0.382 | 0.310 | 0.478 | 0.517 | 0.422 | 0.400 | 0.424 | 0.390 | 0.433 | 0.408 | 0.375 | 0.331 | 0.393 | 0.331 | 0.381 | 0.343 | 7.124e7 | 1.376e16 | 2.278 | 13.843 |
| | ETTm1 (96) | 0.596 | 1.005 | 0.677 | 1.223 | 0.551 | 0.687 | 0.551 | 0.687 | 0.574 | 0.973 | 0.551 | 0.917 | 0.624 | 1.089 | 0.423 | 0.468 | 0.337 | 0.312 | 0.319 | 0.287 | 0.318 | 0.290 | 0.307 | 0.281 | 0.319 | 0.287 | - | - |
| | ETTm1 (192) | 0.640 | 1.107 | 0.741 | 1.423 | 0.567 | 0.728 | 0.567 | 0.728 | 0.617 | 1.067 | 0.629 | 1.093 | 0.670 | 1.197 | 0.447 | 0.487 | 0.358 | 0.341 | 0.356 | 0.343 | 0.345 | 0.325 | 0.355 | 0.324 | 0.389 | 0.394 | 1.050 | 4.497 |
| | ETTm2 (96) | 0.287 | 0.213 | 0.365 | 0.351 | 0.296 | 0.204 | 0.296 | 0.204 | 0.276 | 0.190 | 0.293 | 0.220 | 0.285 | 0.212 | 0.296 | 0.218 | 0.253 | 0.176 | 0.246 | 0.172 | 0.246 | 0.175 | 0.241 | 0.161 | 0.246 | 0.167 | - | - |
| | ETTm2 (192) | 0.322 | 0.268 | 0.432 | 0.482 | 0.341 | 0.273 | 0.341 | 0.273 | 0.317 | 0.250 | 0.344 | 0.311 | 0.324 | 0.268 | 0.347 | 0.306 | 0.292 | 0.231 | 0.278 | 0.217 | 0.277 | 0.215 | 0.274 | 0.204 | 0.291 | 0.226 | 1.806 | 9.276 |
| micro | Male (4) | 15.282 | 1.263e4 | 25.604 | 3.032e4 | 17.842 | 1.414e4 | 17.842 | 1.414e4 | 13.304 | 9.069e3 | 17.054 | 1.231e4 | 16.769 | 1.155e4 | 16.315 | 1.224e4 | 16.121 | 1.293e4 | 14.975 | 1.339e4 | 14.574 | 1.276e4 | 20.197 | 2.279e4 | 32.331 | 9.795e4 | 36.745 | 4.605e4 |
| | Male (8) | 17.280 | 1.426e4 | 27.433 | 3.723e4 | 20.106 | 1.706e4 | 20.106 | 1.706e4 | 16.350 | 1.219e4 | 20.772 | 1.612e4 | 19.314 | 1.429e4 | 17.456 | 1.732e4 | 17.369 | 1.836e4 | 3.051e3 | 8.376e8 | 17.452 | 1.764e4 | 18.820 | 1.854e4 | 17.445 | 1.721e4 | | |
| | Female (4) | 13.405 | 9.310e3 | 21.737 | 2.048e4 | 12.552 | 7.903e3 | 12.552 | 7.903e3 | 10.444 | 6.573e3 | 15.379 | 1.145e4 | 15.540 | 1.491e3 | 28.450 | 4.503e4 | 11.444 | 8.429e3 | 1.406e7 | 1.619e16 | 11.336 | 8.777e3 | 120.134 | 1.157e6 | 13.158 | 8.186e3 | 37.473 | 2.276e4 |
| | Female (8) | 14.123 | 9.312e3 | 21.935 | 2.203e4 | 12.235 | 6.679e3 | 12.235 | 6.679e3 | 11.433 | 6.569e3 | 15.861 | 1.188e4 | 14.534 | 1.064e4 | 12.234 | 6.786e3 | 12.463 | 7.501e3 | 1.915e5 | 1.680e13 | 14.814 | 1.040e4 | 2.768e3 | 3.461e9 | 13.931 | 8.903e3 | | |
| | Donor A (4) | 20.633 | 2.454e4 | 33.762 | 6.267e4 | 18.968 | 1.977e4 | 18.968 | 1.977e4 | 15.124 | 1.461e4 | 21.622 | 3.381e4 | 16.635 | 1.599e4 | 20.269 | 2.275e4 | 34.399 | 6.980e4 | 19.131 | 2.200e4 | 23.232 | 2.772e4 | 21.403 | 2.446e4 | 22.164 | 2.486e4 | 34.559 | 4.047e4 |
| | Donor A (8) | 21.202 | 2.344e4 | 32.012 | 5.255e4 | 18.708 | 1.760e4 | 18.708 | 1.760e4 | 16.953 | 1.670e4 | 22.602 | 3.127e4 | 18.513 | 1.779e4 | 19.290 | 1.931e4 | 19.140 | 2.128e4 | 19.415 | 2.135e4 | 1.145e3 | 7.445e7 | 25.720 | 3.815e4 | 21.405 | 2.554e4 | | |
| | Donor B (4) | 14.202 | 1.441e4 | 21.004 | 3.513e4 | 17.065 | 2.124e4 | 17.065 | 2.124e4 | 12.939 | 1.344e4 | 16.358 | 2.197e4 | 12.118 | 1.162e4 | 19.259 | 2.307e4 | 15.563 | 2.021e4 | 15.208 | 1.827e4 | 15.960 | 2.194e4 | 20.340 | 3.408e4 | 19.002 | 2.558e4 | 38.049 | 5.194e4 |
| | Donor B (8) | 18.515 | 2.313e4 | 26.309 | 4.721e4 | 19.233 | 2.570e4 | 19.233 | 2.570e4 | 16.758 | 1.989e4 | 23.027 | 3.904e4 | 16.884 | 2.069e4 | 21.242 | 2.884e4 | 20.089 | 2.817e4 | 19.967 | 2.901e4 | 19.837 | 2.745e4 | 21.890 | 3.028e4 | 23.713 | 3.457e4 | | |
| epi | Covid (14) | 834.566 | 2.982e7 | 1.222e3 | 4.642e7 | 1.241e3 | 7.355e7 | 1.241e3 | 7.355e7 | 952.373 | 3.320e7 | 844.450 | 2.737e7 | 16629.327 | 1.088e10 | 309.160 | 6.208e5 | 239.591 | 4.463e5 | 250.482 | 4.852e5 | 258.320 | 4.837e5 | 260.914 | 4.764e5 | 311.272 | 6.968e5 | - | - |
| | Covid (30) | 1.736e3 | 1.710e8 | 3.048e3 | 4.914e8 | 1.497e3 | 6.775e7 | 1.497e3 | 6.775e7 | 1.725e3 | 1.431e8 | 2.147e3 | 2.399e8 | 2.811e3 | 2.020e8 | 447.836 | 1.223e6 | 375.605 | 8.041e5 | 335.329 | 6.637e5 | 305.884 | 6.640e5 | 422.727 | 9.982e5 | 581.864 | 4.034e6 | 314.193 | 6.441e5 |
| scRNA | hESC (8) | 3.117 | 17.859 | 4.621 | 37.222 | 2.652 | 10.412 | 2.652 | 10.412 | 2.580 | 10.036 | 3.043 | 14.833 | 2.978 | 11.826 | 2.663 | 10.838 | 2.648 | 11.345 | 2.648 | 11.492 | 2.651 | 11.459 | 2.648 | 10.705 | 2.687 | 10.933 | - | - |
| | hESC (16) | 3.128 | 18.030 | 4.586 | 36.180 | 2.690 | 10.239 | 2.690 | 10.239 | 2.628 | 10.056 | 2.828 | 12.296 | 2.821 | 10.536 | 2.661 | 11.147 | 2.631 | 11.365 | 2.652 | 11.275 | 2.644 | 10.969 | 2.643 | 10.349 | 2.638 | 10.349 | 2.574 | 9.313 |
| | hHep (8) | 2.378 | 11.781 | 3.481 | 23.887 | 2.024 | 7.433 | 2.024 | 7.433 | 1.982 | 7.207 | 2.374 | 9.895 | 2.281 | 8.080 | 2.004 | 7.452 | 2.000 | 7.663 | 2.011 | 7.640 | 2.009 | 7.590 | 2.107 | 7.404 | 2.175 | 7.774 | - | - |
| | hHep (16) | 2.525 | 12.771 | 3.534 | 23.457 | 2.202 | 8.217 | 2.202 | 8.217 | 2.101 | 7.607 | 2.299 | 9.286 | 2.247 | 7.955 | 1.985 | 7.637 | 1.959 | 7.716 | 1.964 | 7.711 | 1.958 | 7.549 | 2.090 | 7.047 | 2.123 | 7.219 | 2.229 | 7.372 |
| | mDC (8) | 1.933 | 6.226 | 2.698 | 12.026 | 1.518 | 3.519 | 1.518 | 3.519 | 1.476 | 3.418 | 1.940 | 6.389 | 1.637 | 4.346 | 1.469 | 3.245 | 1.480 | 3.312 | 1.483 | 3.560 | 1.486 | 3.536 | 1.559 | 3.957 | 1.580 | 4.064 | - | - |
| | mDC (16) | 1.953 | 6.269 | 3.014 | 13.973 | 1.418 | 2.991 | 1.418 | 2.991 | 1.498 | 3.417 | 1.789 | 5.147 | 1.529 | 3.757 | 1.427 | 3.102 | 1.407 | 3.185 | 1.446 | 3.360 | 1.413 | 3.318 | 1.476 | 3.340 | 1.465 | 3.496 | 1.449 | 3.154 |
| | mESC (8) | 0.822 | 1.368 | 1.191 | 2.839 | 0.676 | 0.825 | 0.676 | 0.825 | 0.666 | 0.821 | 0.845 | 1.300 | 0.864 | 1.202 | 0.652 | 0.792 | 0.648 | 0.814 | 0.650 | 0.802 | 0.648 | 0.811 | 0.646 | 0.801 | 0.648 | 0.807 | - | - |
| | mESC (16) | 0.834 | 1.417 | 1.185 | 2.816 | 0.664 | 0.790 | 0.664 | 0.790 | 0.670 | 0.828 | 0.758 | 1.060 | 0.738 | 0.965 | 0.642 | 0.780 | 0.637 | 0.817 | 0.637 | 0.774 | 0.638 | 0.762 | 0.634 | 0.775 | 0.650 | 0.818 | 0.642 | 0.752 |
| | mHSC-E (8) | 3.262 | 17.883 | 4.954 | 39.413 | 2.771 | 10.695 | 2.771 | 10.695 | 2.675 | 10.496 | 2.974 | 13.572 | 2.726 | 9.711 | 2.682 | 10.524 | 2.671 | 10.492 | 2.661 | 10.577 | 2.662 | 10.315 | 2.677 | 10.690 | 2.676 | 10.611 | - | - |
| | mHSC-E (16) | 3.411 | 19.092 | 4.961 | 38.615 | 2.919 | 11.391 | 2.919 | 11.391 | 2.818 | 11.638 | 2.949 | 12.659 | 2.826 | 10.485 | 2.740 | 10.657 | 2.730 | 10.293 | 2.747 | 10.720 | 2.737 | 10.725 | 2.698 | 10.460 | 2.733 | 10.881 | 2.795 | 9.851 |
| | mHSC-GM (8) | 3.362 | 18.373 | 5.057 | 40.229 | 2.816 | 10.799 | 2.816 | 10.799 | 2.719 | 11.127 | 3.121 | 14.811 | 2.769 | 9.737 | 2.693 | 10.876 | 2.683 | 10.517 | 2.676 | 10.743 | 2.684 | 10.669 | 2.702 | 10.860 | 2.681 | 10.806 | - | - |
| | mHSC-GM (16) | 3.397 | 18.754 | 4.872 | 37.419 | 2.827 | 10.535 | 2.827 | 10.535 | 2.788 | 12.029 | 3.029 | 13.303 | 2.782 | 9.832 | 2.722 | 10.905 | 2.699 | 10.629 | 2.707 | 10.622 | 2.711 | 10.264 | 2.692 | 10.883 | 2.711 | 11.195 | 2.694 | 9.292 |
| | mHSC-L (8) | 3.487 | 19.734 | 5.178 | 42.097 | 2.895 | 11.224 | 2.895 | 11.224 | 2.795 | 11.583 | 3.268 | 16.266 | 2.849 | 10.152 | 2.758 | 11.254 | 2.742 | 11.098 | 2.742 | 11.468 | 2.737 | 11.127 | 2.742 | 11.432 | 2.751 | 11.321 | - | - |
| | mHSC-L (16) | 3.461 | 19.292 | 4.898 | 37.801 | 2.859 | 10.587 | 2.859 | 10.587 | 2.824 | 12.207 | 3.098 | 13.762 | 2.841 | 10.084 | 2.727 | 10.912 | 2.715 | 11.277 | 2.717 | 11.182 | 2.720 | 10.893 | 2.703 | 11.059 | 2.771 | 11.437 | 2.748 | 9.575 |

Table 3: Test set MAE and MSE evaluations

8

Table 4: Train set MAPE and sMAPE evaluations

| Group | Dataset (H) | Naive MAPE | Naive sMAPE | RWD MAPE | RWD sMAPE | Average MAPE | Average sMAPE | Window Average MAPE | Window Average sMAPE | ARIMA MAPE | ARIMA sMAPE | ETS MAPE | ETS sMAPE | THETA MAPE | THETA sMAPE | Autoformer MAPE | Autoformer sMAPE | Informer MAPE | Informer sMAPE | N-BEATS MAPE | N-BEATS sMAPE | N-HiTS MAPE | N-HiTS sMAPE | TSMixer MAE | TSMixer sMAPE | Time Mixer MAPE | Time Mixer sMAPE | Time-LLM MAPE | Time-LLM sMAPE |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| dysr | Aizawa (32) | 15.053 | 0.547 | 13.272 | 0.520 | 24.663 | 0.755 | 24.663 | 0.755 | 14.159 | 0.511 | 5.951 | 0.384 | 13.264 | 0.492 | 8.701 | 0.271 | 1.094 | 0.103 | 0.137 | 0.020 | 0.088 | 0.022 | 1.237 | 0.097 | 0.437 | 0.051 | - | - |
| dysr | Aizawa (64) | 16.441 | 0.650 | 20.731 | 0.627 | 23.815 | 0.754 | 23.815 | 0.754 | 15.442 | 0.637 | 9.006 | 0.566 | 18.772 | 0.617 | 21.052 | 0.410 | 2.148 | 0.092 | 0.193 | 0.038 | 0.214 | 0.041 | 2.003 | 0.285 | 1.604 | 0.096 | 17.673 | 0.448 |
| dysr | Blasius (32) | 0.994 | 0.217 | 2.482 | 0.371 | 2.024 | 0.338 | 2.024 | 0.338 | 0.939 | 0.336 | 0.897 | 0.194 | 0.956 | 0.218 | 11.876 | 0.359 | 0.465 | 0.176 | 0.082 | 0.047 | 0.052 | 0.027 | 0.340 | 0.142 | 0.481 | 0.149 | - | - |
| dysr | Blasius (64) | 1.920 | 0.305 | 3.527 | 0.452 | 1.744 | 0.323 | 1.744 | 0.323 | 1.307 | 0.386 | 1.938 | 0.327 | 1.922 | 0.313 | 3.918 | 0.362 | 0.571 | 0.263 | 0.082 | 0.046 | 0.059 | 0.029 | 0.279 | 0.130 | 0.181 | 0.086 | 0.750 | 0.249 |
| dysr | Lorenz (32) | 3.248 | 0.378 | 5.288 | 0.473 | 3.555 | 0.425 | 3.555 | 0.425 | 3.326 | 0.319 | 4.636 | 0.380 | 3.322 | 0.374 | 7.086 | 0.287 | 1.086 | 0.143 | 0.298 | 0.045 | 0.527 | 0.043 | 1.506 | 0.254 | 0.734 | 0.113 | - | - |
| dysr | Lorenz (64) | 2.265 | 0.402 | 3.550 | 0.485 | 3.056 | 0.481 | 3.056 | 0.481 | 2.678 | 0.376 | 3.260 | 0.470 | 2.271 | 0.397 | 3.046 | 0.314 | 1.046 | 0.130 | 0.368 | 0.057 | 0.584 | 0.039 | 1.962 | 0.162 | 0.467 | 0.062 | 0.896 | 0.275 |
| ETT | ETTh1 (96) | 4.306 | 0.492 | 5.611 | 0.538 | 3.598 | 0.513 | 3.598 | 0.513 | 3.622 | 0.480 | 4.271 | 0.484 | 4.445 | 0.490 | 3.517 | 0.443 | 3.180 | 0.424 | 2.737 | 0.397 | 2.753 | 0.397 | 2.907 | 0.400 | 2.908 | 0.401 | 2.834e4 | 0.216 |
| ETT | ETTh1 (192) | 4.747 | 0.516 | 6.379 | 0.553 | 3.529 | 0.524 | 3.529 | 0.524 | 4.091 | 0.504 | 4.707 | 0.510 | 5.044 | 0.515 | 3.342 | 0.462 | 3.251 | 0.443 | 2.812 | 0.431 | 2.783 | 0.426 | 3.006 | 0.421 | 3.138 | 0.426 | - | - |
| ETT | ETTh2 (96) | 2.582 | 0.369 | 3.702 | 0.413 | 2.129 | 0.380 | 2.129 | 0.380 | 2.245 | 0.361 | 2.556 | 0.361 | 2.939 | 0.374 | 2.115 | 0.380 | 2.164 | 0.363 | 1.789 | 0.310 | 1.917 | 0.318 | 1.641 | 0.300 | 1.956 | 0.325 | 8.672e4 | 0.126 |
| ETT | ETTh2 (192) | 2.732 | 0.400 | 3.735 | 0.433 | 2.540 | 0.410 | 2.540 | 0.410 | 2.416 | 0.389 | 2.720 | 0.396 | 3.095 | 0.407 | 2.818 | 0.411 | 1.267 | 0.235 | 1.501 | 0.358 | 1.514 | 0.360 | 1.681 | 0.305 | 5.575e8 | 1.000 | - | - |
| ETT | ETTm1 (96) | 3.808 | 0.455 | 5.313 | 0.494 | 2.692 | 0.482 | 2.692 | 0.482 | 3.803 | 0.445 | 4.238 | 0.460 | 3.780 | 0.447 | 2.229 | 0.391 | 1.724 | 0.320 | 1.925 | 0.284 | 2.075 | 0.356 | 1.865 | 0.342 | 1.732 | 0.308 | 10024.971 | 0.195 |
| ETT | ETTm1 (192) | 4.130 | 0.493 | 5.251 | 0.522 | 2.736 | 0.500 | 2.736 | 0.500 | 4.121 | 0.483 | 4.578 | 0.498 | 4.033 | 0.486 | 2.490 | 0.428 | 2.090 | 0.377 | 1.925 | 0.352 | 2.075 | 0.356 | 1.865 | 0.342 | 1.633 | 0.310 | - | - |
| ETT | ETTm2 (96) | 2.058 | 0.339 | 2.637 | 0.376 | 1.891 | 0.355 | 1.891 | 0.355 | 1.978 | 0.334 | 2.038 | 0.339 | 2.068 | 0.341 | 1.750 | 0.329 | 1.179 | 0.246 | 1.413 | 0.280 | 1.452 | 0.283 | 1.443 | 0.288 | 1.389 | 0.282 | 2.929e5 | 0.184 |
| ETT | ETTm2 (192) | 2.328 | 0.365 | 2.913 | 0.395 | 2.184 | 0.379 | 2.184 | 0.379 | 2.239 | 0.359 | 2.309 | 0.364 | 2.388 | 0.370 | 1.936 | 0.344 | 1.731 | 0.325 | 1.574 | 0.302 | 1.575 | 0.305 | 1.545 | 0.305 | 1.668 | 0.318 | - | - |
| micro | Male (4) | 1.675 | 0.243 | 3.375 | 0.338 | 1.891 | 0.315 | 1.891 | 0.315 | 1.193 | 0.268 | 2.155 | 0.802 | 2.222 | 0.334 | 1.714 | 0.811 | 1.248 | 0.844 | 1.147 | 0.830 | 1.147 | 0.833 | 0.971 | 0.789 | 1.036 | 0.793 | 4.16e6 | 0.658 |
| micro | Male (8) | 1.937 | 0.265 | 3.707 | 0.360 | 2.258 | 0.379 | 2.258 | 0.379 | 1.442 | 0.341 | 2.307 | 0.809 | 2.666 | 0.412 | 1.211 | 0.834 | 1.097 | 0.841 | 432.355 | 0.983 | 1.059 | 0.834 | 1.747 | 0.821 | 1.454 | 0.823 | - | - |
| micro | Female (4) | 1.248 | 0.119 | 2.060 | 0.161 | 1.371 | 0.160 | 1.371 | 0.160 | 0.915 | 0.121 | 1.427 | 0.871 | 1.435 | 0.171 | 2.935 | 0.947 | 1.026 | 0.886 | 1.835e6 | 1.000 | 1.010 | 0.886 | 0.725 | 0.860 | 0.803 | 0.862 | 2.995e6 | 0.641 |
| micro | Female (8) | 1.561 | 0.134 | 2.633 | 0.183 | 1.651 | 0.207 | 1.651 | 0.207 | 1.102 | 0.158 | 1.458 | 0.880 | 1.857 | 0.225 | 1.333 | 0.886 | 1.149 | 0.889 | 1.657e4 | 0.999 | 1.446 | 0.896 | 0.788 | 0.838 | 0.888 | 0.870 | - | - |
| micro | Donor A (4) | 1.836 | 0.160 | 3.293 | 0.221 | 1.979 | 0.215 | 1.979 | 0.215 | 1.182 | 0.173 | 1.760 | 0.879 | 2.005 | 0.248 | 1.664 | 0.898 | 1.000 | 1.000 | 1.232 | 0.902 | 2.242 | 0.888 | 0.944 | 0.863 | 1.243 | 0.875 | 2.144e6 | 0.794 |
| micro | Donor A (8) | 1.984 | 0.166 | 3.286 | 0.222 | 2.183 | 0.267 | 2.183 | 0.267 | 1.431 | 0.215 | 1.910 | 0.882 | 2.319 | 0.304 | 1.880 | 0.890 | 1.383 | 0.900 | 1.277 | 0.898 | 169.700 | 0.981 | 0.973 | 0.877 | 0.935 | 0.876 | - | - |
| micro | Donor B (4) | 1.688 | 0.197 | 3.384 | 0.272 | 1.743 | 0.250 | 1.743 | 0.250 | 1.095 | 0.206 | 1.300 | 0.826 | 1.710 | 0.295 | 1.564 | 0.839 | 1.158 | 0.866 | 1.236 | 0.863 | 1.125 | 0.870 | 1.219 | 0.856 | 1.162 | 0.837 | 2.033e6 | 0.702 |
| micro | Donor B (8) | 1.786 | 0.212 | 2.793 | 0.292 | 1.743 | 0.308 | 1.743 | 0.308 | 1.272 | 0.261 | 1.498 | 0.838 | 1.809 | 0.351 | 1.505 | 0.853 | 1.145 | 0.870 | 1.122 | 0.875 | 1.142 | 0.871 | 1.154 | 0.869 | 0.889 | 0.836 | - | - |
| epi | Covid (14) | 0.109 | 0.050 | 0.154 | 0.073 | 0.107 | 0.050 | 0.107 | 0.050 | 0.107 | 0.054 | 7.902 | 0.111 | 0.127 | 0.059 | 0.105 | 0.074 | 0.094 | 0.070 | 0.103 | 0.075 | 0.106 | 0.076 | 0.108 | 0.075 | 0.066 | 0.055 | 0.134 | 0.062 |
| epi | Covid (30) | 0.115 | 0.054 | 0.158 | 0.073 | 0.129 | 0.058 | 0.129 | 0.058 | 0.132 | 0.072 | 0.592 | 0.115 | 0.125 | 0.057 | 0.154 | 0.095 | 0.072 | 0.061 | 0.004 | 0.026 | 0.110 | 0.080 | 0.154 | 0.092 | 0.078 | 0.062 | - | - |
| scRNA | hESC (8) | 0.788 | 0.398 | 1.172 | 0.497 | 0.599 | 0.394 | 0.599 | 0.394 | 0.582 | 0.393 | 0.635 | 0.428 | 0.706 | 0.431 | 0.603 | 0.447 | 0.631 | 0.447 | 0.624 | 0.451 | 0.629 | 0.438 | 0.612 | 0.441 | 0.587 | 0.440 | 1.070e6 | 0.316 |
| scRNA | hESC (16) | 0.800 | 0.406 | 1.147 | 0.500 | 0.592 | 0.403 | 0.592 | 0.403 | 0.581 | 0.407 | 0.606 | 0.419 | 0.650 | 0.425 | 0.625 | 0.440 | 0.618 | 0.451 | 0.608 | 0.441 | 0.625 | 0.438 | 0.584 | 0.432 | 0.595 | 0.438 | - | - |
| scRNA | hHep (8) | 1.600 | 0.424 | 2.257 | 0.519 | 1.417 | 0.453 | 1.417 | 0.453 | 1.387 | 0.456 | 1.358 | 0.514 | 1.518 | 0.494 | 1.342 | 0.512 | 1.393 | 0.506 | 1.386 | 0.507 | 1.379 | 0.505 | 1.343 | 0.505 | 1.032 | 0.477 | 1.742e6 | 0.368 |
| scRNA | hHep (16) | 1.592 | 0.454 | 2.212 | 0.551 | 1.404 | 0.477 | 1.404 | 0.477 | 1.389 | 0.488 | 1.354 | 0.518 | 1.443 | 0.501 | 1.385 | 0.509 | 1.386 | 0.510 | 1.398 | 0.514 | 1.359 | 0.515 | 1.190 | 0.477 | 1.324 | 0.500 | - | - |
| scRNA | mDC (8) | 1.405e11 | 0.427 | 1.736e11 | 0.516 | 1.195e11 | 0.355 | 1.195e11 | 0.355 | 1.191e11 | 0.361 | 1.614e11 | 0.344 | 1.722e11 | 0.403 | 1.077e11 | 0.351 | 1.132e11 | 0.353 | 1.092e11 | 0.351 | 1.190e11 | 0.354 | 1.662e10 | 0.290 | 1.188e10 | 0.327 | 4.936e5 | 0.297 |
| scRNA | mDC (16) | 1.302e11 | 0.463 | 1.419e11 | 0.533 | 1.326e11 | 0.356 | 1.326e11 | 0.356 | 1.315e11 | 0.371 | 1.586e11 | 0.347 | 8.086e10 | 0.378 | 1.185e11 | 0.350 | 1.278e11 | 0.355 | 1.208e11 | 0.370 | 1.310e11 | 0.360 | 9.096e10 | 0.359 | 2.471e10 | 0.338 | - | - |
| scRNA | mESC (8) | 6.605e5 | 0.335 | 9.866e5 | 0.428 | 7.824e5 | 0.314 | 7.824e5 | 0.314 | 7.273e5 | 0.319 | 6.718e5 | 0.333 | 9.174e5 | 0.350 | 7.278e5 | 0.332 | 7.104e5 | 0.339 | 7.218e5 | 0.336 | 7.120e5 | 0.341 | 7.071e5 | 0.319 | 7.775e5 | 0.335 | 6.083e5 | 0.254 |
| scRNA | mESC (16) | 6.545e5 | 0.339 | 1.052e6 | 0.425 | 8.161e5 | 0.319 | 8.161e5 | 0.319 | 7.793e5 | 0.330 | 7.945e5 | 0.326 | 9.630e5 | 0.341 | 7.542e5 | 0.331 | 6.758e5 | 0.347 | 7.162e5 | 0.335 | 7.690e5 | 0.329 | 8.097e5 | 0.333 | 6.568e5 | 0.327 | - | - |
| scRNA | mHSC-E (8) | 1.254 | 0.480 | 1.949 | 0.611 | 1.059 | 0.417 | 1.059 | 0.417 | 0.934 | 0.421 | 1.007 | 0.412 | 1.121 | 0.449 | 0.928 | 0.439 | 0.869 | 0.429 | 0.872 | 0.434 | 0.930 | 0.428 | 0.958 | 0.436 | 0.979 | 0.434 | 7.285e5 | 0.367 |
| scRNA | mHSC-E (16) | 1.281 | 0.481 | 1.916 | 0.609 | 1.087 | 0.423 | 1.087 | 0.423 | 0.929 | 0.435 | 1.053 | 0.417 | 1.140 | 0.445 | 0.908 | 0.467 | 0.924 | 0.431 | 0.902 | 0.434 | 0.895 | 0.438 | 0.936 | 0.422 | 0.947 | 0.426 | - | - |
| scRNA | mHSC-GM (8) | 1.268 | 0.516 | 1.985 | 0.639 | 1.014 | 0.447 | 1.014 | 0.447 | 0.865 | 0.462 | 0.986 | 0.440 | 1.173 | 0.485 | 0.829 | 0.461 | 0.835 | 0.461 | 0.777 | 0.464 | 0.808 | 0.455 | 0.855 | 0.447 | 0.910 | 0.466 | 7.214e5 | 0.367 |
| scRNA | mHSC-GM (16) | 1.275 | 0.521 | 1.909 | 0.637 | 1.017 | 0.444 | 1.017 | 0.444 | 0.827 | 0.479 | 0.992 | 0.440 | 1.109 | 0.473 | 0.831 | 0.464 | 0.813 | 0.460 | 0.833 | 0.459 | 0.864 | 0.455 | 0.823 | 0.445 | 0.871 | 0.464 | - | - |
| scRNA | mHSC-L (8) | 1.273 | 0.525 | 1.997 | 0.653 | 1.010 | 0.459 | 1.010 | 0.459 | 0.859 | 0.477 | 0.972 | 0.451 | 1.178 | 0.502 | 0.817 | 0.479 | 0.787 | 0.477 | 0.747 | 0.480 | 0.781 | 0.472 | 0.881 | 0.483 | 0.904 | 0.479 | 8.997e5 | 0.375 |
| scRNA | mHSC-L (16) | 1.280 | 0.524 | 1.916 | 0.645 | 0.999 | 0.454 | 0.999 | 0.454 | 0.819 | 0.494 | 0.981 | 0.451 | 1.105 | 0.485 | 0.812 | 0.474 | 0.766 | 0.477 | 0.768 | 0.477 | 0.786 | 0.472 | 0.821 | 0.466 | 0.833 | 0.467 | - | - |

Table 4: Train set MAPE and sMAPE evaluations

Table 5 — Test set MAPE and sMAPE evaluations

| | Dataset (H) | Naive MAE | Naive MSE | RWD MAE | RWD MSE | Average MAE | Average MSE | Window Average MAE | Window Average MSE | ARIMA MAE | ARIMA MSE | ETS MAE | ETS MSE | THETA MAE | THETA MSE | Autoformer MAE | Autoformer MSE | Informer MAE | Informer MSE | N-BEATS MAE | N-BEATS MSE | N-HITS MAE | N-HITS MSE | TSMixer MAE | TSMixer MSE | Time Mixer MAE | Time Mixer MSE | Time-LLM MAE | Time-LLM MSE |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| dysl | Aizawa (32) | 3.931 | 0.573 | 7.657 | 0.555 | 5.433 | 0.766 | 5.433 | 0.766 | 3.646 | 0.539 | 5.542 | 0.381 | 5.137 | 0.518 | 1.848 | 0.270 | 0.557 | 0.126 | 0.088 | 0.028 | 0.106 | 0.027 | 0.581 | 0.106 | 0.210 | 0.061 | - | 0.465 |
| | Aizawa (64) | 4.884 | 0.654 | 7.201 | 0.624 | 5.259 | 0.737 | 5.259 | 0.737 | 4.327 | 0.642 | 8.902 | 0.574 | 4.853 | 0.626 | 2.142 | 0.466 | 0.382 | 0.107 | 0.162 | 0.047 | 0.150 | 0.046 | 1.070 | 0.282 | 0.394 | 0.099 | 11.872 | - |
| | Blasius (32) | 1.215 | 0.242 | 2.981 | 0.405 | 2.212 | 0.371 | 2.212 | 0.371 | 1.072 | 0.340 | 1.108 | 0.220 | 1.178 | 0.248 | 9.628 | 0.337 | 0.406 | 0.172 | 0.062 | 0.032 | 0.108 | 0.052 | 0.361 | 0.167 | 0.665 | 0.157 | - | 0.207 |
| | Blasius (64) | 2.190 | 0.326 | 3.811 | 0.483 | 1.750 | 0.335 | 1.750 | 0.335 | 1.392 | 0.400 | 2.203 | 0.348 | 2.203 | 0.334 | 4.301 | 0.361 | 0.549 | 0.241 | 0.069 | 0.044 | 0.046 | 0.025 | 0.313 | 0.156 | 0.148 | 0.074 | 0.690 | - |
| | Lorenz (32) | 30.601 | 0.439 | 39.277 | 0.505 | 35.978 | 0.508 | 35.978 | 0.508 | 9.447 | 0.360 | 7.200 | 0.403 | 30.822 | 0.424 | 1.162 | 0.332 | 0.621 | 0.164 | 0.158 | 0.052 | 0.120 | 0.043 | 1.066 | 0.266 | 0.341 | 0.108 | - | 0.333 |
| | Lorenz (64) | 9.482 | 0.474 | 4.379 | 0.476 | 25.110 | 0.504 | 25.110 | 0.504 | 12.382 | 0.429 | 11.901 | 0.492 | 9.569 | 0.468 | 0.824 | 0.304 | 0.451 | 0.197 | 0.263 | 0.122 | 0.154 | 0.085 | 0.517 | 0.209 | 0.234 | 0.111 | 1.587 | - |
| ETT | ETTh1 (96) | 3.973 | 0.476 | 4.442 | 0.508 | 3.005 | 0.511 | 3.005 | 0.511 | 3.211 | 0.487 | 3.829 | 0.470 | 3.947 | 0.472 | 2.530 | 0.390 | 2.286 | 0.384 | 2.014 | 0.359 | 1.988 | 0.368 | 1.940 | 0.346 | 2.163 | 0.349 | - | 0.173 |
| | ETTh1 (192) | 4.025 | 0.486 | 4.829 | 0.522 | 3.026 | 0.534 | 3.026 | 0.534 | 2.990 | 0.511 | 3.897 | 0.483 | 3.987 | 0.482 | 2.079 | 0.405 | 2.624 | 0.392 | 2.243 | 0.393 | 2.295 | 0.374 | 2.449 | 0.362 | 2.581 | 0.372 | 0.521 | - |
| | ETTh2 (96) | 2.176 | 0.333 | 3.105 | 0.387 | 1.960 | 0.332 | 1.960 | 0.332 | 2.004 | 0.324 | 2.203 | 0.332 | 2.461 | 0.348 | 1.478 | 0.310 | 1.346 | 0.301 | 1.069 | 0.270 | 1.129 | 0.268 | 1.164 | 0.277 | 1.074 | 0.269 | - | 0.351 |
| | ETTh2 (192) | 3.226 | 0.357 | 4.363 | 0.430 | 3.054 | 0.363 | 3.054 | 0.363 | 3.116 | 0.340 | 3.174 | 0.351 | 3.448 | 0.373 | 1.346 | 0.355 | 1.565 | 0.341 | 1.270 | 0.318 | 1.391 | 0.322 | 1.266 | 0.316 | 3.416e8 | 1.000 | 1.852e5 | - |
| | ETTm1 (96) | 3.610 | 0.442 | 4.130 | 0.471 | 2.769 | 0.492 | 2.769 | 0.492 | 3.388 | 0.435 | 3.835 | 0.447 | 3.078 | 0.435 | 2.363 | 0.395 | 1.983 | 0.332 | 1.857 | 0.319 | 1.901 | 0.316 | 1.897 | 0.308 | 2.012 | 0.313 | - | 0.160 |
| | ETTm1 (192) | 3.827 | 0.463 | 4.367 | 0.497 | 2.774 | 0.505 | 2.774 | 0.505 | 3.598 | 0.457 | 4.070 | 0.469 | 3.645 | 0.464 | 2.356 | 0.418 | 2.062 | 0.345 | 2.103 | 0.345 | 2.239 | 0.329 | 2.141 | 0.347 | 2.411 | 0.357 | 1.778e3 | - |
| | ETTm2 (96) | 1.524 | 0.244 | 1.896 | 0.279 | 1.314 | 0.266 | 1.314 | 0.266 | 1.404 | 0.246 | 1.516 | 0.242 | 1.463 | 0.253 | 1.303 | 0.249 | 1.187 | 0.223 | 1.153 | 0.210 | 1.142 | 0.208 | 1.084 | 0.206 | 1.052 | 0.210 | - | 0.341 |
| | ETTm2 (192) | 1.665 | 0.265 | 2.155 | 0.309 | 1.414 | 0.288 | 1.414 | 0.288 | 1.562 | 0.268 | 1.662 | 0.266 | 1.718 | 0.279 | 1.450 | 0.266 | 1.237 | 0.236 | 1.155 | 0.224 | 1.217 | 0.224 | 1.120 | 0.222 | 1.136 | 0.232 | 4.320e5 | - |
| micro | Male (4) | 1.141 | 0.224 | 2.422 | 0.326 | 1.445 | 0.306 | 1.445 | 0.306 | 0.891 | 0.265 | 1.781 | 0.780 | 1.631 | 0.332 | 1.304 | 0.787 | 0.999 | 0.816 | 0.858 | 0.801 | 0.865 | 0.803 | 2.350 | 0.798 | 6.431 | 0.802 | - | 0.636 |
| | Male (8) | 1.245 | 0.248 | 2.221 | 0.350 | 1.645 | 0.382 | 1.645 | 0.382 | 1.067 | 0.345 | 1.668 | 0.792 | 1.983 | 0.417 | 0.954 | 0.813 | 0.915 | 0.818 | 293.033 | 0.982 | 0.910 | 0.816 | 1.269 | 0.803 | 1.117 | 0.803 | 2.239e7 | - |
| | Female (4) | 1.363 | 0.149 | 2.324 | 0.203 | 1.314 | 0.201 | 1.314 | 0.201 | 0.974 | 0.146 | 1.621 | 0.864 | 1.860 | 0.229 | 3.201 | 0.944 | 1.025 | 0.878 | 1.707e6 | 1.000 | 1.015 | 0.879 | 11.477 | 0.939 | 1.442 | 0.877 | - | 0.626 |
| | Female (8) | 1.545 | 0.154 | 2.710 | 0.199 | 1.263 | 0.247 | 1.263 | 0.247 | 1.124 | 0.172 | 1.477 | 0.869 | 1.944 | 0.286 | 1.241 | 0.874 | 1.092 | 0.880 | 1.540e4 | 0.899 | 1.610 | 0.889 | 206.614 | 0.974 | 1.354 | 0.880 | 3.018e6 | - |
| | Donor A (4) | 1.741 | 0.164 | 3.281 | 0.234 | 1.927 | 0.223 | 1.927 | 0.223 | 1.268 | 0.183 | 1.519 | 0.873 | 1.962 | 0.265 | 1.487 | 0.895 | 1.000 | 1.000 | 1.186 | 0.897 | 2.214 | 0.885 | 2.075 | 0.891 | 2.175 | 0.899 | - | 0.782 |
| | Donor A (8) | 1.944 | 0.175 | 3.311 | 0.241 | 2.037 | 0.278 | 2.037 | 0.278 | 1.451 | 0.227 | 1.767 | 0.876 | 2.154 | 0.321 | 1.652 | 0.884 | 1.263 | 0.896 | 1.289 | 0.897 | 161.310 | 0.982 | 2.639 | 0.900 | 1.552 | 0.892 | 2.281e6 | - |
| | Donor B (4) | 1.206 | 0.139 | 1.960 | 0.190 | 1.127 | 0.183 | 1.127 | 0.183 | 0.996 | 0.161 | 1.029 | 0.872 | 1.292 | 0.373 | 0.949 | 0.852 | 0.811 | 0.878 | 0.836 | 0.873 | 0.795 | 0.882 | 1.086 | 0.883 | 1.148 | 0.876 | - | 0.755 |
| | Donor B (8) | 1.464 | 0.161 | 2.404 | 0.213 | 1.289 | 0.229 | 1.289 | 0.229 | 1.283 | 0.198 | 1.326 | 0.885 | 1.734 | 0.429 | 1.230 | 0.880 | 1.021 | 0.899 | 0.941 | 0.895 | 0.948 | 0.894 | 1.112 | 0.896 | 1.395 | 0.887 | 5.037e6 | - |
| epi | Covid (14) | 0.159 | 0.079 | 0.277 | 0.142 | 0.191 | 0.080 | 0.191 | 0.080 | 0.161 | 0.081 | 7.334 | 0.100 | 0.176 | 0.095 | 0.100 | 0.049 | 0.077 | 0.038 | 0.081 | 0.041 | 0.084 | 0.042 | 0.086 | 0.042 | 0.096 | 0.047 | - | 0.055 |
| | Covid (30) | 0.248 | 0.089 | 0.487 | 0.142 | 0.289 | 0.103 | 0.289 | 0.103 | 0.261 | 0.111 | 0.704 | 0.123 | 0.353 | 0.125 | 0.141 | 0.069 | 0.124 | 0.065 | 0.118 | 0.058 | 0.098 | 0.050 | 0.132 | 0.062 | 0.177 | 0.068 | 0.111 | - |
| scRNA | hESC (8) | 0.906 | 0.458 | 1.356 | 0.569 | 0.712 | 0.476 | 0.712 | 0.476 | 0.690 | 0.472 | 0.752 | 0.486 | 0.825 | 0.511 | 0.668 | 0.497 | 0.698 | 0.499 | 0.696 | 0.508 | 0.699 | 0.506 | 0.679 | 0.489 | 0.678 | 0.491 | - | 0.351 |
| | hESC (16) | 0.904 | 0.467 | 1.330 | 0.591 | 0.717 | 0.484 | 0.717 | 0.484 | 0.700 | 0.484 | 0.733 | 0.481 | 0.764 | 0.509 | 0.683 | 0.489 | 0.669 | 0.505 | 0.679 | 0.498 | 0.685 | 0.485 | 0.654 | 0.482 | 0.644 | 0.484 | 1447129.000 | - |
| | hHep (8) | 2.188 | 0.410 | 2.801 | 0.501 | 1.815 | 0.485 | 1.815 | 0.485 | 1.748 | 0.481 | 1.692 | 0.563 | 1.779 | 0.535 | 2.376 | 0.629 | 2.450 | 0.630 | 2.543 | 0.627 | 2.477 | 0.624 | 2.252 | 0.612 | 2.207 | 0.610 | - | 0.444 |
| | hHep (16) | 2.217 | 0.431 | 3.049 | 0.531 | 1.658 | 0.529 | 1.658 | 0.529 | 1.694 | 0.526 | 1.777 | 0.563 | 1.757 | 0.555 | 2.495 | 0.630 | 2.545 | 0.637 | 2.499 | 0.642 | 2.310 | 0.642 | 2.257 | 0.606 | 2.262 | 0.606 | 1975036.000 | - |
| | mDC (8) | 3.299e9 | 0.510 | 3.529e9 | 0.585 | 1.997e9 | 0.426 | 1.997e9 | 0.426 | 1.982e9 | 0.425 | 3.310e9 | 0.406 | 2.931e9 | 0.481 | 3.793e9 | 0.423 | 3.923e9 | 0.429 | 4.892e9 | 0.434 | 3.617e9 | 0.427 | 3.896e9 | 0.454 | 4.207e9 | 0.462 | - | 0.374 |
| | mDC (16) | 3.530e9 | 0.545 | 4.480e9 | 0.657 | 2.408e9 | 0.402 | 2.408e9 | 0.402 | 2.385e9 | 0.436 | 3.305e9 | 0.396 | 2.590e9 | 0.494 | 4.012e9 | 0.412 | 4.243e9 | 0.419 | 3.564e9 | 0.450 | 4.777e9 | 0.424 | 2.907e9 | 0.433 | 4.143e9 | 0.437 | 5.775e5 | - |
| | mESC (8) | 2.975e5 | 0.354 | 4.840e5 | 0.444 | 3.818e5 | 0.349 | 3.818e5 | 0.349 | 3.664e5 | 0.352 | 5.847e5 | 0.396 | 4.805e5 | 0.391 | 3.871e5 | 0.370 | 3.778e5 | 0.383 | 3.964e5 | 0.376 | 3.704e5 | 0.384 | 4.001e5 | 0.375 | 4.157e5 | 0.376 | - | 0.274 |
| | mESC (16) | 5.913e5 | 0.354 | 8.742e5 | 0.441 | 4.170e5 | 0.352 | 4.170e5 | 0.352 | 3.775e5 | 0.363 | 5.212e5 | 0.369 | 5.389e5 | 0.380 | 2.883e5 | 0.372 | 1.631e5 | 0.391 | 2.551e5 | 0.375 | 2.977e5 | 0.367 | 3.025e5 | 0.371 | 2.898e5 | 0.383 | 7.421e5 | - |
| | mHSC-E (8) | 1.248 | 0.474 | 1.953 | 0.599 | 1.016 | 0.411 | 1.016 | 0.411 | 0.905 | 0.414 | 1.001 | 0.406 | 1.119 | 0.443 | 0.870 | 0.391 | 0.824 | 0.387 | 0.814 | 0.388 | 0.879 | 0.384 | 0.887 | 0.394 | 0.905 | 0.392 | - | 0.356 |
| | mHSC-E (16) | 1.286 | 0.486 | 1.915 | 0.608 | 1.062 | 0.421 | 1.062 | 0.421 | 0.896 | 0.431 | 1.039 | 0.412 | 1.071 | 0.439 | 0.857 | 0.394 | 0.894 | 0.389 | 0.873 | 0.391 | 0.850 | 0.394 | 0.914 | 0.391 | 0.946 | 0.395 | 5.578e5 | - |
| | mHSC-GM (8) | 1.328 | 0.498 | 2.071 | 0.630 | 1.089 | 0.426 | 1.089 | 0.426 | 0.904 | 0.440 | 1.066 | 0.421 | 1.232 | 0.471 | 0.847 | 0.412 | 0.873 | 0.406 | 0.800 | 0.409 | 0.830 | 0.412 | 0.937 | 0.415 | 0.920 | 0.414 | - | 0.349 |
| | mHSC-GM (16) | 1.308 | 0.502 | 1.959 | 0.625 | 1.074 | 0.424 | 1.074 | 0.424 | 0.850 | 0.456 | 1.061 | 0.421 | 1.169 | 0.458 | 0.848 | 0.413 | 0.836 | 0.410 | 0.865 | 0.408 | 0.890 | 0.407 | 0.883 | 0.413 | 0.881 | 0.420 | 4.820e5 | - |
| | mHSC-L (8) | 1.276 | 0.508 | 1.965 | 0.631 | 1.012 | 0.439 | 1.012 | 0.439 | 0.860 | 0.453 | 0.991 | 0.433 | 1.210 | 0.479 | 0.805 | 0.424 | 0.778 | 0.422 | 0.733 | 0.423 | 0.774 | 0.420 | 0.856 | 0.431 | 0.881 | 0.428 | - | 0.356 |
| | mHSC-L (16) | 1.294 | 0.499 | 1.885 | 0.613 | 0.997 | 0.434 | 0.997 | 0.434 | 0.815 | 0.464 | 0.991 | 0.431 | 1.133 | 0.464 | 0.801 | 0.421 | 0.755 | 0.421 | 0.763 | 0.421 | 0.788 | 0.418 | 0.834 | 0.423 | 0.901 | 0.428 | 6.198e5 | - |

Table 5: Test set MAPE and sMAPE evaluations