

SPECIALIZED FOUNDATION MODELS STRUGGLE TO BEAT SUPERVISED BASELINES

Anonymous authors

Paper under double-blind review

ABSTRACT

Following its success for vision and text, the “foundation model” (FM) paradigm—pretraining large models on massive data, then fine-tuning on target tasks—has rapidly expanded to domains in the sciences, engineering, healthcare, and beyond. Has this achieved what the original FMs accomplished, i.e. the supplanting of traditional supervised learning in their domains? To answer we look at three modalities—genomics, satellite data, and time series—with multiple recent FMs and compare them to a standard supervised learning workflow: model development, hyperparameter tuning, and training, all using only data from the target task. Across those three specialized domains, we find that it is consistently possible to train simple supervised models—no more complicated than a lightly modified wide ResNet or UNet—that match or even outperform the latest foundation models. Our work demonstrates that the benefits of large-scale pretraining have yet to be realized in many specialized areas, reinforces the need to compare new FMs to strong, well-tuned baselines, and introduces two new, easy-to-use, open-source, and automated workflows for doing so.

1 INTRODUCTION

Recent years have witnessed a shift towards large-scale pretraining across domains like computer vision and natural language processing. This workflow generally consists of two stages: pretraining on vast amounts of domain-specific data to capture general knowledge followed by fine-tuning on target tasks (Radford & Narasimhan, 2018). This pretrain-then-finetune paradigm has been tremendously successful, enabling foundation models (Bommasani et al., 2021) to consistently outcompete traditional supervised learning methods on a wide variety of downstream tasks in the vision and language domains (Dosovitskiy et al., 2021; Liu et al., 2021; Devlin et al., 2019).

Driven by this success, the foundation model approach has been adapted to various *specialized* domains, which we define to be ML application areas—e.g. genomics, satellite imaging, and time series—whose data modalities lie outside those of classical AI tasks, i.e. natural images and text. These domains have seen the introduction of many new FMs claiming to leverage large, domain-specific pretraining datasets to achieve breakthrough performance on downstream tasks (Dalla-Torre et al., 2023; Nguyen et al., 2024; Zhou et al., 2023b; Avsec et al., 2021; Ji et al., 2021; Fuller et al., 2023; Cong et al., 2022; Mendieta et al., 2023). This claim underlies our study’s motivating question:

Do these new specialized FMs outperform traditional supervised learning applied to the same tasks?

Answering this question is critical because supervised workflows are usually much less expensive to implement and deploy, but if FMs do dominate them then FMs have the potential to fundamentally transform these domains, as we have seen with language and vision processing in the past decade. However, despite ongoing efforts to promote their fair and comprehensive evaluation (Liang et al., 2022; Bommasani & Liang, 2021), many new FMs have not been adequately compared to simpler, often more efficient baselines. Indeed, we found that many works only benchmark their proposed models against other FMs, essentially creating a comparison echo chamber (Fuller et al., 2023; Mendieta et al., 2023; Nguyen et al., 2024; Zhou et al., 2023b).

We answer our motivating question by considering a reasonably representative set of three specialized domains—chosen according to the presence of multiple FMs and a standard set of evaluation tasks—and comparing their performance on those tasks with that of a traditional supervised learning

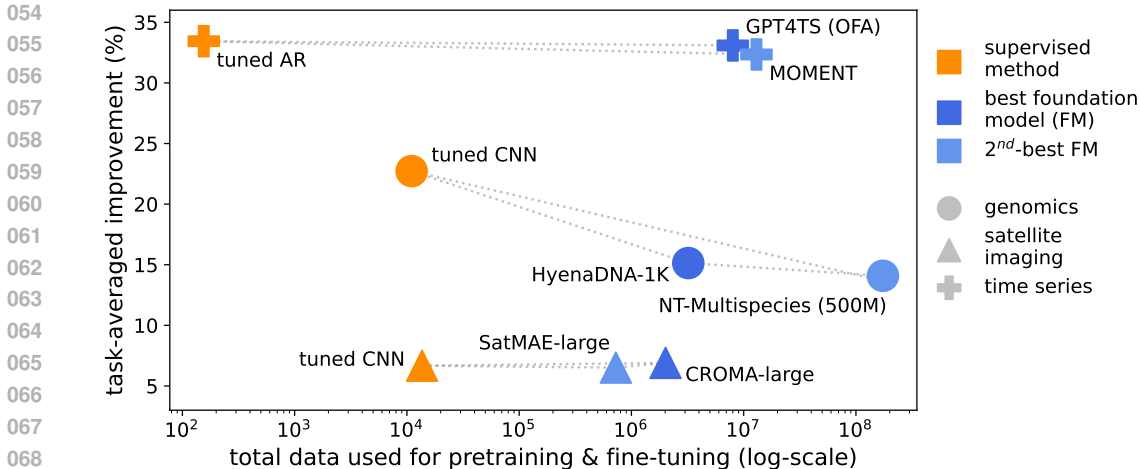


Figure 1: Across three domains, tuned supervised baselines match or outperform the best specialized FMs while using two-to-five orders of magnitude less data. We plot total pretraining and fine-tuning data (in kilobase-pairs for genomics, images for satellite, and unique series for time series) vs. the mean percentage improvement (in MCC, accuracy / mAP, and RMSE, respectively) over an untuned domain-specific baseline (respectively a 1D Wide ResNet, a 2D Wide ResNet, and Auto-ARIMA).

workflow. As depicted in Figure 2, the latter is a model development, hyperparameter tuning, and training process in which all steps use only data from the target task, in contrast to the FM workflow, which uses vast amounts of pretraining data. By leveraging model selection tools ranging from classical information criteria to cutting-edge architecture search, we build automated pipelines that efficiently develop and train strong supervised models on over fifty tasks across three distinct domains.

Our high-level result is negative: we find that, despite being pretrained on massive datasets, specialized FMs struggle and very often fail to outperform models trained exclusively on downstream task data with traditional supervised learning (c.f. Figure 1). Specifically, we show that lightly adapted convolutional neural network (CNN) architectures such as wide ResNet and UNet attain state-of-the-art on the Nucleotide Transformer benchmark in genomics and match the latest pretrained satellite FMs on downstream classification. Furthermore, we show that tuned linear auto-regression (AR) outperforms all open-source time series FMs on a standard suite of eight forecasting tasks and approaches the performance of two other closed-source models that only evaluate on subsets of them.

These results demonstrate that genomics, satellite imaging, and time series have not yet had their “BERT moment” (Devlin et al., 2019), i.e. these domains have not yet pretrained FMs that dominate traditional supervised approaches. This is despite the fact that all them have BERT-scale FMs—hundreds of millions of parameters or larger—and the fact that many of them are already witnessing a shift towards not comparing with supervised approaches, as was seen in natural language processing (NLP) post-BERT. More broadly, since these domains are among the most high-profile areas with specialized FMs, our results challenge the prevailing assumption that pretraining them has led to superior performance. They also reinforce the need for robust and well-tuned baselines, with surprising findings such as (a) simply tuning kernel sizes and dilation rates in standard CNN backbones dominates a genomics classification benchmark and (b) rescuing the century-old AR forecaster from obsolescence is as easy as considering lookback parameters larger than five and training on a GPU. To facilitate ongoing research in these domains and others, we will make code associated with both our CNN-tuning pipeline (DASHA) and our AR-on-GPU workflow (Auto-AR) publicly available.

2 RELATED WORK

Foundation models have been trained in numerous specialized domains beyond vision and text, including genomics (Ji et al., 2021), satellite imaging (Cong et al., 2022), time series (Goswami et al., 2024), weather (Bodnar et al., 2024), differential equation solving (Sun et al., 2024), web traffic (Zhao et al., 2023), and beyond. To get a representative sense of their success, we focus on

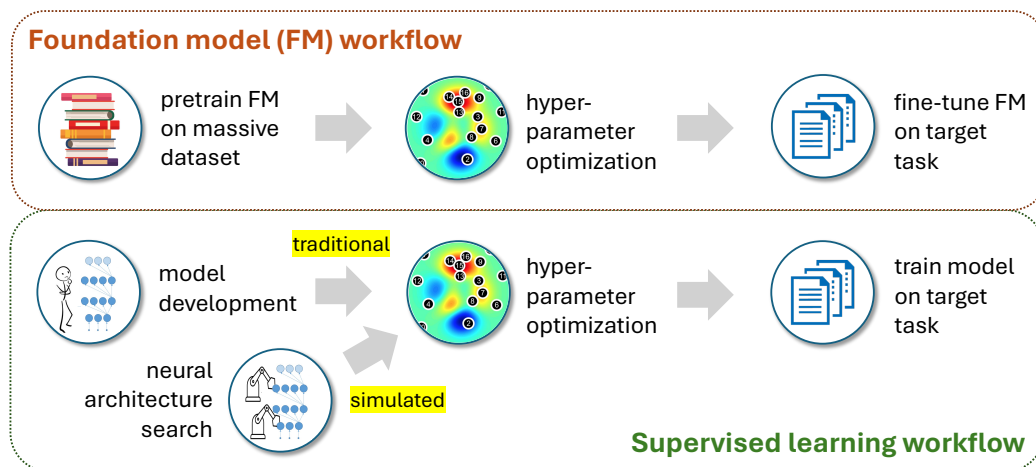


Figure 2: Our goal is to compare the pretrain-then-fine-tune paradigm (top) with a standard supervised workflow (bottom) on the tasks on which specialized FMs are evaluated. While for time series we go through a traditional process of developing and tuning a supervised model, this manual approach does not scale to many domains; as a result, in Section 3.1 we develop a way to simulate it using architecture search. Note that FM fine-tuning hyperparameters are not always tuned in practice, but we assume their creators make a best-effort attempt to present their own method in the best light.

domains that combine the following properties: (a) multiple BERT-scale FMs, (b) a standard suite of evaluation tasks, and (c) significant applied interest. This restriction naturally suggest looking at three domains that all have at least five FMs evaluated on at least nine tasks: genomics (which has some of the largest-available non-text FMs (Dalla-Torre et al., 2023)), satellite imaging (which has a large ongoing benchmarking effort (Lacoste et al., 2024)), and time series (which has already seen significant industry interest (Cohen et al., 2024)). The remainder of this section examines how different learning workflows approach problems in these domains.

2.1 SPECIALIZED FOUNDATION MODELS

Our three target domains have more than twenty FMs between them, many developed via the “lift-and-shift” approach—borrowing terminology from Rolf et al. (2024)—in which techniques from core AI areas such as vision and language processing are applied with modest tailoring to specialized domains. In particular, many methods are built around out-of-domain models such as BERT, Swin, and Hyena (Ji et al., 2021; Mendieta et al., 2023; Nguyen et al., 2024), with adaptations such as specialized tokenizations, embeddings, and model modifications for handling domain-specific considerations such long-range dependencies (Dalla-Torre et al., 2023; Zhou et al., 2023b; Das et al., 2023; Cohen et al., 2024) or multispectral data (Cong et al., 2022).

While the “lift-and-shift” approach can often be useful or at least a good starting point, its widespread use underlines the need for strong *in-domain* baselines to make sure that the combination of out-of-domain tooling and massive pretraining data is actually helpful. Such comparisons are not always conducted, e.g. the satellite FM SatMAE (Cong et al., 2022) is compared to ImageNet-initialized and randomly initialized ResNet-50 (He et al., 2015), while most of the time series FMs we consider only do a full comparison to one linear baseline, DLinear (Zeng et al., 2023). While this can sometimes be justified—e.g. in the case of NLP post-BERT—our results suggest that for now, specialized FMs should still compare to in-domain supervised model development.

2.2 SPECIALIZED BASELINES

Both of the automated supervised learning pipelines we develop are heavily influenced by successful in-domain model development. In particular, the NAS-based pipeline we use to achieve our results in genomics and satellite imaging is inspired by the success of the human-driven specification of kernel sizes and dilation rates in successful architectures like TCN (Lea et al., 2016) and ConvNeXt (Liu

162 *et al.*, 2022). At the same time, for time series our approach is based upon a well-tuned GPU
163 implementation of perhaps the most basic forecasting model, AR.
164

165 2.3 AUTOML FOR SPECIALIZED DOMAINS 166

167 While often evaluated on domains such as vision, automated techniques have long been used in
168 specialized domains as well. An important example is Auto-ARIMA (Hyndman & Khandakar,
169 2008) for time series, although it has been found to underperform on the specific suite of tasks we
170 consider (Challu *et al.*, 2022). However, to avoid requiring significant expertise in any one domain,
171 we also make use of AutoML methods developed specifically for diverse tasks (Roberts *et al.*,
172 2021b; Shen *et al.*, 2023), in particular the NAS method DASH (Shen *et al.*, 2022) that can discover
173 good kernel sizes and dilation rates for a CNN backbone faster than it can be trained from scratch.
174

175 3 METHODOLOGY 176

177 Recall that our goal is to conduct a robust comparison between traditional supervised learning and
178 specialized FMs; the natural way to do this is to take existing benchmarks used to evaluate FMs in
179 our three target domains and run a typical supervised workflow on the same tasks. As depicted in
180 Figure 2, this pipeline involves three steps: (1) model development, (2) hyperparameter tuning, and
181 (3) training. The first stage involves using both reasoning and trial-and-error to find a good architec-
182 ture to tune and train on the data; for example, Lea *et al.* (2016) developed the temporal convolutional
183 network (TCN) architecture with a multi-layer dilation rate pattern specifically suited to sequential
184 data, while Liu *et al.* (2022) designed the breakthrough ConvNeXt architecture by methodically
185 exploring ways to make CNNs more like Transformers without introducing attention. The second
186 stage (hyperparameter tuning) can also be done via human-driven iteration, but there exist effective
187 automated procedures for it as well (Li *et al.*, 2020). Lastly, the third step of the pipeline involves
188 simply training the selected model with the selected configuration on the data of the target task.

189 While it standard to automate the last two steps of the procedure, model development is typically
190 done by hand and so is difficult to do for fifty tasks across three domains. As a result, we settle for
191 *approximating* the traditional supervised learning workflow by simulating the model development
192 component using neural architecture search. To ensure fair comparison and reduce computational
193 costs, we restrict ourselves to low-fidelity NAS methods that return an architecture in less time than
194 it takes to train it. The results we obtain using NAS can therefore be viewed as *lower bounds* on
195 the performance of supervised learning, as the model development might be significantly improved
196 using less-heuristic or human-driven architecture design.

197 In the remainder of this section we detail how we handle the different steps of the supervised learning
198 pipeline. Note that our NAS-dependent supervised workflow (DASHA)—which we cover in the
199 first part of this section—yields our main results for genomics and satellite imaging but *not* for time
200 series; in that domain we find its performance to be less competitive. There we instead focus on an
201 even simpler approach based on linear auto-regression, whose model development and tuning we
202 describe in the second subsection.

203 3.1 DASHA: SIMULATING THE SUPERVISED WORKFLOW USING NAS 204

205 To simulate model development we need a search space over architectures that is (a) efficient,
206 (b) flexible, and (c) applicable to the types of high-dimensional unstructured data that arise in
207 domains targeted by specialized FMs; these requirements make CNN-based search spaces a nat-
208 ural choice. In particular, inspired by the success of hand-tuned kernel sizes and dilation rates
209 in traditional model development (Lea *et al.*, 2016; Bai *et al.*, 2018; Liu *et al.*, 2022), we apply
210 DASH (Shen *et al.*, 2022), a NAS method that starts with an existing CNN backbone—e.g. a wide
211 ResNet (Zagoruyko & Komodakis, 2017)—and uses the weight-sharing heuristic (Liu *et al.*, 2018)
212 to determine the right kernel size and dilation rate to use at each convolutional layer. DASH has been
213 successfully used in AutoML competitions (Roberts *et al.*, 2021a) and to advance the state-of-the-art
214 on NAS benchmarks (Tu *et al.*, 2022), making it likely to be useful beyond the domains we consider.

215 As described in Algorithm 1, we augment the existing DASH approach in two ways: (1) trying
more than one CNN backbone (e.g. both wide ResNet and U-Net (Ronneberger *et al.*, 2015)) and

Algorithm 2: Pseudocode for the DASHA workflow. Starting with a set of backbone CNNs, we use DASH (Shen et al., 2022) to set the right kernel size and dilation rate for each of its convolutional layers and then use ASHA (Li et al., 2020) to configure a training routine for the resulting architecture. Lastly, we pick the best backbone using validation data and train it.

Input: target task dataset D , candidate CNN backbone architectures A

for CNN backbone $a \in A$ **do**

 // set a kernel size and dilation rate for each layer of a

$\text{arch}_a \leftarrow \text{DASH}(D, a)$

 // tune hyperparameters for the discovered architecture arch_a

$\text{config}_a, \text{val_score}_a \leftarrow \text{ASHA}(D, \text{arch}_a)$

 // train the architecture with the highest validation score

$a \leftarrow \arg \max_{a \in A} \text{val_score}_a$

Output: $\text{train}(D, \text{arch}_a, \text{config}_a)$

(2) using the well-known hyperparameter tuner ASHA (Li et al., 2020) to configure architecture-specific training settings. This combination gives our workflow its name. Following the NAS and hyperparameter tuning stages, we train the discovered architecture with the selected configuration on the target data. Further details, including the resources given to the three steps of the pipeline and the exact search spaces used by DASH and ASHA, are provided in Appendix A.1. Note that, while our focus is on *data*-efficient baselines, we do ensure that the entire workflow is never substantially more computationally expensive than fine-tuning an FM.

3.2 AUTO-AR: MAKING A BASELINE STRONGER BY MAKING IT SIMPLER

While DASHA can be applied to forecasting tasks, it is not competitive with state-of-the-art time series FMs. At the same time, the field of time series forecasting has long employed automated workflows, notably the Auto-ARIMA approach of Hyndman & Khandakar (2008) that uses statistical tests and information criteria to tune ARIMA’s lookback and differencing parameters. Auto-ARIMA was evaluated on the time series tasks we consider by Challu et al. (2022), who found that it performed poorly compared to deep learning approaches. However, their implementation does not make use of multi-channel data and tunes up to a lookback window of at most five, which is much less data than used by time series FMs. While tuning ARIMA with larger lookback parameters is computationally costly, we find the following simplified tuning pipeline to be effective:

1. use the KPSS test (Kwiatkowski et al., 1992) to decide whether to take first differences
2. use the Bayesian Information Criterion to select the maximum lookback parameter of the auto-regressive (AR) component of ARIMA, ignoring the moving average (MA) part
3. maximize the multi-channel likelihood of AR with the chosen differencing and lookback

By dropping the MA component of the model and running the procedure on GPU, we are able to tune the lookback windows up the maximum allowable length (usually 512); we find that longer lookbacks are critical for performance. Note that this is just a tuned version of the classic AR model.

4 EMPIRICAL RESULTS

We now present the results of applying the automated pipelines described in the previous section to our three target domains. For each domain, we provide a brief justification of the specific FMs and evaluation tasks that we consider, followed by details on how we apply our workflows; further information can be found in Appendix A and Appendix B. As there are too many separate results to present outside the appendix, in this section we mainly present aggregate statistics that summarize our findings for each domain, with detailed results relegated to Appendix B. The domains have different performance metrics, but they can all be aggregated via the following quantities: **average score**, **average rank**, and **mean / median percentage improvement over a baseline**. For each domain, we define a domain-specific baseline and measure the improvement of FMs and our approach relative to it. This standardizes comparisons across tasks of varying scales.

Model	Model Size	Pretraining Base-Pairs	Avg. MCC \uparrow	Avg. Rank \downarrow	Mean %Imp. \uparrow	Median %Imp. \uparrow
Foundation Models						
NT-1000G (500M)	500M	20.5T	0.625	8.42	2.15	-6.41
NT-1000G (2.5B)	2.5B	20.5T	0.656	5.47	7.15	-0.83
DNABERT-2	117M	32.5B	0.680	5.39	10.40	-0.81
HyenaDNA-1K	436K	3.2B	0.708	5.47	15.14	-2.63
HyenaDNA-32K	1.6M	3.2B	0.630	8.17	2.55	-11.12
Enformer	–	4B	0.568	9.75	-8.58	-19.26
NT-Multispecies (2.5B)	2.5B	174B	0.697	3.19	13.66	1.59
NT-Multispecies (500M)	500M	174B	0.700	2.97	14.09	2.28
Supervised Methods						
Wide ResNet	2.0M	0	0.657	5.69	0.00	0.00
UNet	4.5M	0	0.614	7.39	-6.94	-1.45
DASHA (our workflow)	10.5M	0	0.752	4.08	22.58	2.66

Table 1: Aggregate performance on genomics tasks, showing that our supervised workflow (DASHA) attains state-of-the-art on the NT benchmark, outperforming all FMs according to most measures while using no pretraining data and oftentimes many fewer parameters. For Mean/Median %Imp. we report percentage improvement over a vanilla 1D Wide ResNet baseline, and for DASHA the model size refers to the largest configuration across tasks. “–” indicates unknown quantities.

4.1 GENOMICS

We begin our investigation in the genomics domain, which has witnessed the development of numerous FMs, including the early Enformer [Avsec et al. \(2021\)](#), the DNABERT series ([Ji et al., 2021](#); [Zhou et al., 2023b](#)), the HyenaDNA family ([Nguyen et al., 2024](#)), and the NT family ([Dalla-Torre et al., 2023](#)); the latter includes models with up to 2.5B parameters. To evaluate them, we consider the Nucleotide Transformer (NT) benchmark of [Dalla-Torre et al. \(2023\)](#), which contains eighteen tasks in three main categories: regulatory elements, RNA production, and histone modification. We use this benchmark because of its diversity and because it has been evaluated on by all of the aforementioned FMs, allowing us to include eight of them in the comparison. Our numbers for these models are taken from [Dalla-Torre et al. \(2023, Supplementary Table 6\)](#); we use Matthew’s Correlation Coefficient (MCC) as the main metric for evaluation.

4.1.1 BASELINES

CNNs have long been used for genomics tasks ([Avsec et al., 2020](#); [Zhou & Troyanskaya, 2015](#)) and so constitute natural supervised baselines; in particular we include 1D variants of Wide ResNet (WRN) and UNet, which we find perform better than some domain-specific CNNs. We use these same two backbones as the candidate CNNs tuned and selected from by our DASHA workflow.

4.1.2 RESULTS

Our genomics results are displayed in [Table 1](#), which shows that our supervised workflow (DASHA) outperforms all FMs across all aggregate metrics except average rank, where it lags behind the two NT-Multispecies models. Notably, according to the last column those are also the only two FMs that even improve over Wide ResNet on the typical task in the NT benchmark. As discussed in [Appendix B](#), our strong performance is driven in large part by outstanding performance on the histone modification tasks (c.f. [Table 9](#)). The more detailed results also highlight the importance of considering diverse baselines, with Wide ResNet usually being the selected architecture but UNet performing significantly better for promoter and splice site classification tasks. Overall, DASHA arguably sets a new state-of-the-art on the NT benchmark and certainly demonstrates that supervised methods remain quite competitive in genomics, despite the availability of massive pretraining datasets.

Model	Model Size	Pretraining Images	Average Score \uparrow	Average Rank \downarrow	Mean %Imp. \uparrow	Median %Imp. \uparrow
Foundation Models						
CROMA-base	90.6M	2M	77.39	4.33	5.85	4.22
CROMA-large	312M	2M	78.03	3.33	6.90	6.09
SatMAE-base	85.6M	700K	76.99	6.22	5.27	3.59
SatMAE-large	303M	700K	77.75	4.5	6.52	4.62
GFM	86.8M	1.3M	77.18	5.56	5.77	4.08
SwinT-base	86.8M	14M	76.69	5.28	4.86	1.43
Supervised Methods						
ResNet50	23.5M	0	73.76	8.34	0.30	00.07
Wide ResNet	17.2M	0	73.97	8.22	0.00	0.00
UNet	17.3M	0	75.73	5.89	3.01	1.07
DASHA (our workflow)	32.4M	0	77.85	3.33	6.67	5.16

Table 2: Aggregate performance on satellite imaging tasks, demonstrating that a supervised learning workflow (DASHA) can match the performance of state-of-the-art specialized FMs, all while using no pretraining data and having two-to-ten times fewer parameters. For Mean/Median %Imp. we report percentage improvement over a vanilla Wide ResNet, and for DASHA the model size refers to the largest configuration across tasks.

4.2 SATELLITE IMAGING

While they do not get as large as those in genomics, numerous BERT-scale FMs have also been introduced for satellite imaging, including SeCo (Manas et al., 2021), the SatMAE family (Cong et al., 2022), the CROMA family (Fuller et al., 2023), GFM (Mendieta et al., 2023), Scale-MAE (Reed et al., 2023), Satlas (Bastani et al., 2023), Prithvi (Jakubik et al., 2023), and SkySense (Guo et al., 2024). Because our evaluation includes GeoBench (Lacoste et al., 2024), a recently introduced satellite benchmark that has not been considered by many of these FMs, we must obtain all the results using our own fine-tuning; therefore we only consider a restricted subset of top-performing, open-source, and compatibly-formatted models. In all cases we use the fine-tuning workflow suggested by the authors of each FM plus some automated hyperparameter tuning. We take our tasks mainly from GeoBench’s five classification tasks and then add four additional tasks—BigEarthNet (Sumbul et al., 2019), EuroSAT (Helber et al., 2019), Canadian Cropland (Jacques et al., 2023), and fMoW-Sentinel (Cong et al., 2022)—that are commonly used to evaluate other FMs.¹ As we focus on classification—sometimes with multiple labels—we report top-1 accuracy or mAP as appropriate.

4.2.1 BASELINES

Since satellite imaging closely resembles RGB imaging, it is common to “lift-and-shift” vision models to this domain (Rolf et al., 2024). As a result we use several CNN backbones as additional baselines, and use wide ResNet as the candidate architecture for our DASHA workflow. Note specifically that because we treat ResNet50 architectures as one of our supervised baselines, it is trained with random initialization. Lastly, we also consider the performance of fine-tuning the vision FM SwinT-base (Liu et al., 2021), which is pretrained on ImageNet.

4.2.2 RESULTS

As shown by Table 2, our supervised workflow attains the best or second-best performance across all aggregate metrics and is only ever slightly outperformed by CROMA-large. Notably, unlike in genomics, the FMs here consistently outperform CNN backbones, likely because the associated papers compare to them as baselines. However, the frequently superior performance of DASHA here suggests that domain-aware model development would yield effective supervised models in this field. Another contrast with genomics is that the larger versions of different FMs consistently attain superior performance here, suggesting they are making at least somewhat effective use of the pretraining data. However, the fact that this improvement can also be attained by DASHA, which

¹In Appendix B we report results when excluding tasks where missing channels may affect performance.

Model	Model Size	Pretraining Series	Full setting (8 datasets / 32 tasks)				Partial setting (6 datasets / 24 tasks)			
			Avg. RMSE ↓	Avg. Rank ↓	Mean %Imp. ↑	Median %Imp. ↑	Avg. RMSE ↓	Avg. Rank ↓	Mean %Imp. ↑	Median %Imp. ↑
Foundation Models										
GPT4TS (OFA)	60M	8M	0.659	3.22	33.09	31.38	0.540	6.23	27.30	16.60
LLM4TS	60M	8M	–	–	–	–	0.526	2.46	29.25	20.96
MOMENT	385M	13M	0.689	2.81	32.35	26.99	0.534	4.56	28.07	18.66
TEMPO (Zero Shot)	–	–	0.733	6.06	27.63	24.51	0.579	9.42	22.69	15.72
TimesFM (Zero Shot)	200M	5M	0.720	4.81	28.94	27.67	0.569	7.98	23.61	17.33
Moirai (Zero Shot)	14M	6M	–	–	–	–	0.566	7.46	24.53	18.52
Toto (Zero Shot)	103M	–	–	–	–	–	0.505	4.38	32.40	31.35
Supervised Methods										
Auto-ARIMA	10	0	1.104	7.94	0.00	0.00	0.806	11.88	0.00	0.00
AR	513	0	0.680	3.66	32.23	30.06	0.540	5.79	27.24	17.28
DLinear	700k	0	0.680	4.42	31.27	30.50	0.551	7.29	25.59	16.94
Auto-AR	513	0	0.660	3.08	33.45	32.20	0.534	4.94	28.12	19.63
DASHA	480K	0	–	–	–	–	0.549	5.62	25.94	16.78

Table 3: Aggregate performance on time series tasks across both six tasks (ETT, Weather, & ECL) and eight tasks (those plus ILI & Traffic). The latter evaluation demonstrates that simply tuning a classical AR model is competitive with state-of-the-art FMs while using no pretraining data and tens of thousands of times fewer parameters. For Mean/Median %Imp. we report percentage improvement over Auto-ARIMA, and for both Auto-AR and DASHA the model size refers to the largest configuration across tasks. “–” indicates unknown quantities.

uses no pretraining data and produces a model that is ten times smaller, shows that there remains significant room for improvement.

4.3 TIME SERIES

The last domain we investigate is time series, where multiple FMs have been introduced, including those that use the standard pretrain-then-fine-tune workflow (GPT4TS (OFA) (Zhou et al., 2023a), LLM4TS (Chang et al., 2023), MOMENT (Goswami et al., 2024), and Time-LLM (Jin et al., 2024))² and others that evaluate in a zero-shot (ZS) regime (TEMPO (Cao et al., 2024), TimesFM (Das et al., 2024), Moirai (Woo et al., 2024), and Toto (Cohen et al., 2024)). As we are comparing to supervised baselines, our evaluation of ZS models will of course be in a less challenging setting than the one they report numbers for. We study the performance of these models and our baselines on the problem of long-horizon forecasting, for which there exists a standard set of datasets summarized in Goswami et al. (2024, Table 11).³ Our main results are on eight of these datasets, but we also report aggregate performance on a subset of six in order to include FMs that do not report numbers on one or more of the other two. Note that each dataset consists of four tasks corresponding to four different time horizons, so in total this yields thirty-two tasks. Lastly, we compute aggregate metrics using RMSE, instead of MSE, so that performance scales linearly with prediction error; this choice has no effect on average rank.

4.3.1 BASELINES

To baseline these FMs we use mainly linear forecasting methods, including the classical (untuned) linear auto-regression (AR), the automated workhorse method Auto-ARIMA (Hyndman & Khandakar, 2008), the more recent DLinear (Zeng et al., 2023), and our own workflow Auto-AR described in Section 3.2. Lastly, we also evaluate our other automated workflow, DASHA, on six of the tasks.

4.3.2 RESULTS

Table 3 shows that on the full eight-dataset evaluation our Auto-AR workflow always attains the best or second best performance across all aggregate metrics considered, and in particular attains the best average improvement over Auto-ARIMA. Notably, the two closest methods in this subset are *not* zero-shot, which is perhaps not surprising given the extra data. However, it does reinforce

²We do not compare to Time-LLM because they report MAE whereas most models focus on MSE.

³We do not consider one of them, Exchange, because most FMs do not report performance on it.

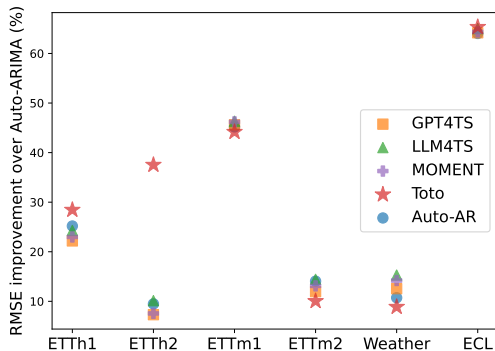


Figure 3: Scatterplot depicting performance (averaged across time horizons) of the best time series methods on the six-dataset subset from Table 3. This visually demonstrates that the recent Toto method wins according to aggregate metrics mainly due to its dominant performance on ETTh2; elsewhere it is middle-of-the-pack.

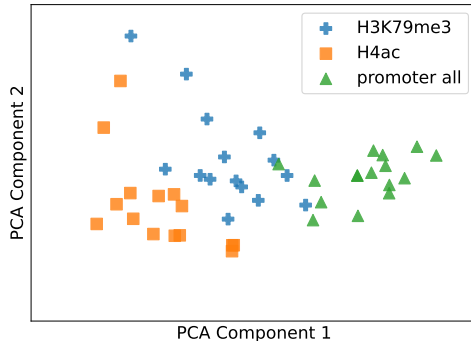


Figure 4: PCA visualization of the architectures discovered for three different tasks when the DASHA workflow is run multiples times. The spatial clustering across tasks demonstrates the within-task consistency of the workflow’s architecture search component and the utility of using diverse models as baselines.

the intuition that settings with high data availability should prefer supervised methods, including simple ones like AR. Notably, even our untuned implementation of AR that uses no differencing and a large lookback window is quite effective, doing better than ZS FMs across all aggregate metrics and even MOMENT on some of them.

In the more limited six-dataset evaluation we consider three additional FMs, two of which—LLM4TS and Toto—are closed-source and thus challenging to evaluate on the other tasks. While our supervised Auto-AR workflow is reasonably close to LLM4TS across most metrics, Toto dominates all metrics except average rank; impressively, it does this despite being zero-shot. However, a look at Figure 3 reveals that Toto is not truly dominant on these six tasks, with its aggregate metrics being strongly influenced by its dramatically better performance on ETTh2, a dataset where all other FMs struggle to do more than 10% better than Auto-ARIMA. Thus, while its ZS performance is quite good, it is unclear whether its domination of aggregate metrics would continue with additional tasks.

5 DISCUSSION

At a high level, our results show that the foundation models in these three domains have not yet surpassed supervised learning, and thus more broadly that the latter remains a strong baseline for specialized FMs. This is a surprising and consequential finding due the paradigm’s popularity and the data and compute costs associated with large-scale pretraining. In this section we discuss lessons and implications for the development of machine learning in these and other application areas.

5.1 THE IMPORTANCE OF DIVERSE, WELL-TUNED, AND DOMAIN-SPECIFIC BASELINES

The main lesson of our work is to select a diverse array of baselines, drawing from both “lift-and-shift” and domain-specific approaches, and then to carefully tune them. For example, in genomics the vanilla wide ResNet baseline does remarkably well, with the majority of FMs doing worse than even this “lift-and-shift” baseline on the typical task in the NT benchmark. While satellite FMs do outperform such baselines, lightly modifying these CNNs via different kernel sizes and dilation rates was enough to match state-of-the-art models there as well. Lastly, our time series results demonstrate in dramatic fashion the need to carefully tune domain-specific approaches, as we show that simply allowing the classical AR forecaster to make use of long lookback windows and GPU-based optimization leads better forecasting than all open-source FMs.

5.2 COMPUTATIONAL EFFICIENCY CONSIDERATIONS

While not the main focus of our work, it is nevertheless worth highlighting that any performance gains of using FMs must be balanced against their additional cost. In addition to the extensive GPU-hours used for pretraining, the resulting models are often much bigger and so lead to much more costly inference. Indeed, apart from the special case of HyenaDNA, the CNN architectures discovered and trained using our DASHA workflow are typically over ten times smaller than FMs in the case of genomics and three to ten times smaller in the case of satellite imaging. Moreover, for time series our Auto-AR approach is quick-to-train and yields simple models with less than a 1K parameters—over ten-thousand times smaller than any FM—while attaining performance that on some datasets is competitive even with closed-source models. In aggregate, these examples further demonstrate the efficiency of supervised approaches and the resulting high performance bar that FMs need to clear before they can be deemed useful.

5.3 THE POWER OF TUNING KERNEL SIZES AND DILATION RATES

Our results in the first two domains, genomics and satellite imaging, are driven by the DASHA workflow, whose crucial component is the tuning of kernel sizes and dilation rate in CNN backbones such as wide ResNet. Its success demonstrates that the procedure is an effective surrogate for human-driven model development, enabling the automated discovery of the types of diverse and in some sense domain-specific baselines stressed in Section 5.1. To understand this further, we study whether the architecture search component selects different kernel sizes and dilation rates for different tasks, and whether it does so in a consistent manner. Specifically, we run DASHA on three of the smaller datasets in the NT benchmark with fifteen different random seeds, construct eighteen-dimensional vectors of the discovered kernel sizes and dilation rates assigned to each of the nine layers, and project these to two dimensions using principal component analysis (PCA). The result in Figure 4 reveals that the architectures are clustered by task, demonstrating that the procedure selects different but consistent-within-task kernel parameters. This visualization suggests that architecture search is a useful surrogate for model development, and consequently that the DASHA workflow may also be useful for automating similar studies and baselining FMs in other domains with high-dimensional, unstructured data.

5.4 THE SURPRISING EFFECTIVENESS OF LINEAR AUTO-REGRESSION

Perhaps our most surprising finding is the competitiveness of linear auto-regression (AR), a very old method, on long-horizon forecasting. It is likely that the lack of comparison with this baseline was driven by existing evaluations (e.g. by Challu et al. (2022)) of Auto-ARIMA (Hyndman & Khandakar, 2008), which is *perceived* to be a stronger baseline because it both combines AR with another model (MA) and tunes the lookback and differencing parameters. However, in most Auto-ARIMA packages the default maximum lookback is around five, whereas we often found much (hundred-fold) larger settings to work best. Since these implementations are also generally too slow to support such long lookbacks, the possibility of expanding the hyperparameter space was more likely to be ignored. By implementing an efficient tuning procedure over a larger space of lookback parameters, our Auto-AR workflow comprises a significant contribution to forecasting baselines.

6 CONCLUSION

We conduct a thorough investigation to evaluate whether the cost of training specialized foundation models across three major domains are justified by their superior performance relative to traditional supervised learning. Our results demonstrate that FMs in these domains have not yet surpassed supervised workflows and are often outperformed by fairly simple methods, including lightly modified CNN backbones (in genomics and satellite imaging) and classical linear forecasters (for time series). As part of our study, we introduce two automated workflows—DASHA for simulating in-domain model development of CNNs and Auto-AR for tuning linear auto-regression on GPUs—that we believe will be useful tools for evaluating future work in these and other areas. The code for these pipelines and to reproduce our results will be made publicly available.

REFERENCES

- 540
541
542 Žiga Avsec, Melanie Weilert, Avanti Shrikumar, Sabrina Krueger, Amr Alexandari, Khyati Dalal,
543 Robin Fropf, Charles McAnany, Julien Gagneur, Anshul Kundaje, and Julia Zeitlinger. Base-
544 resolution models of transcription factor binding reveal soft motif syntax. *bioRxiv*, 2020. doi:
545 10.1101/737981. URL [https://www.biorxiv.org/content/early/2020/07/19/](https://www.biorxiv.org/content/early/2020/07/19/737981)
546 [737981](https://www.biorxiv.org/content/early/2020/07/19/737981). 6
- 547 Žiga Avsec, Vikram Agarwal, Daniel Visentin, Joseph R. Ledsam, Agnieszka Grabska-Barwinska,
548 Kyle R. Taylor, Yannis Assael, John Jumper, Pushmeet Kohli, and David R. Kelley. Effective gene
549 expression prediction from sequence by integrating long-range interactions. *Nature Methods*, 18
550 (10):1196–1203, Oct 2021. doi: 10.1038/s41592-021-01252-x. 1, 6
- 551 Shaojie Bai, J. Zico Kolter, and Vladlen Koltun. An empirical evaluation of generic convolutional
552 and recurrent networks for sequence modeling. arXiv, 2018. 4
- 553
554 Favyen Bastani, Piper Wolters, Ritwik Gupta, Joe Ferdinando, and Aniruddha Kembhavi. Satlaspre-
555 train: A large-scale dataset for remote sensing image understanding. In *IEEE/CVF International*
556 *Conference on Computer Vision, ICCV 2023, Paris, France, October 1-6, 2023*, pp. 16726–16736.
557 IEEE, 2023. doi: 10.1109/ICCV51070.2023.01538. URL [https://doi.org/10.1109/](https://doi.org/10.1109/ICCV51070.2023.01538)
558 [ICCV51070.2023.01538](https://doi.org/10.1109/ICCV51070.2023.01538). 7
- 559 Cristian Bodnar, Wessel P Bruinsma, Ana Lucic, Megan Stanley, Johannes Brandstetter, Patrick
560 Garvan, Maik Riechert, Jonathan Weyn, Haiyu Dong, Anna Vaughan, et al. Aurora: A foundation
561 model of the atmosphere. *arXiv preprint arXiv:2405.13063*, 2024. 2
- 562
563 Rishi Bommasani and Percy Liang. Reflections on foundation models, 2021. URL [https://](https://hai.stanford.edu/news/reflections-foundation-models)
564 hai.stanford.edu/news/reflections-foundation-models. 1
- 565
566 Rishi Bommasani, Drew A. Hudson, Ehsan Adeli, Russ Altman, Simran Arora, Sydney von Arx,
567 Michael S. Bernstein, Jeannette Bohg, Antoine Bosselut, Emma Brunskill, Erik Brynjolfsson,
568 S. Buch, Dallas Card, Rodrigo Castellon, Niladri S. Chatterji, Annie S. Chen, Kathleen A. Creel,
569 Jared Davis, Dora Demszky, Chris Donahue, Moussa Doumbouya, Esin Durmus, Stefano Ermon,
570 John Etchemendy, Kawin Ethayarajh, Li Fei-Fei, Chelsea Finn, Trevor Gale, Lauren E. Gillespie,
571 Karan Goel, Noah D. Goodman, Shelby Grossman, Neel Guha, Tatsunori Hashimoto, Peter Hen-
572 derson, John Hewitt, Daniel E. Ho, Jenny Hong, Kyle Hsu, Jing Huang, Thomas F. Icard, Saahil
573 Jain, Dan Jurafsky, Pratyusha Kalluri, Siddharth Karamcheti, Geoff Keeling, Fereshte Khani,
574 O. Khattab, Pang Wei Koh, Mark S. Krass, Ranjay Krishna, Rohith Kuditipudi, Ananya Kumar,
575 Faisal Ladhak, Mina Lee, Tony Lee, Jure Leskovec, Isabelle Levent, Xiang Lisa Li, Xuechen
576 Li, Tengyu Ma, Ali Malik, Christopher D. Manning, Suvir P. Mirchandani, Eric Mitchell, Zanele
577 Munyikwa, Suraj Nair, Avaniika Narayan, Deepak Narayanan, Benjamin Newman, Allen Nie,
578 Juan Carlos Niebles, Hamed Nilforoshan, J. F. Nyarko, Giray Ogut, Laurel Orr, Isabel Papadimi-
579 triou, Joon Sung Park, Chris Piech, Eva Portelance, Christopher Potts, Aditi Raghunathan, Robert
580 Reich, Hongyu Ren, Frieda Rong, Yusuf H. Roohani, Camilo Ruiz, Jack Ryan, Christopher
581 R’e, Dorsa Sadigh, Shiori Sagawa, Keshav Santhanam, Andy Shih, Krishna Parasuram Sriniva-
582 san, Alex Tamkin, Rohan Taori, Armin W. Thomas, Florian Tramèr, Rose E. Wang, William
583 Wang, Bohan Wu, Jiajun Wu, Yuhuai Wu, Sang Michael Xie, Michihiro Yasunaga, Jiaxuan You,
584 Matei A. Zaharia, Michael Zhang, Tianyi Zhang, Xikun Zhang, Yuhui Zhang, Lucia Zheng, Kait-
585 lyn Zhou, and Percy Liang. On the opportunities and risks of foundation models. *ArXiv*, 2021.
586 URL <https://crfm.stanford.edu/assets/report.pdf>. 1
- 587
588 George.E.P. Box and Gwilym M. Jenkins. *Time Series Analysis: Forecasting and Control*. Holden-
589 Day, 1976. 19
- 590 Defu Cao, Furong Jia, Serkan O Arik, Tomas Pfister, Yixiang Zheng, Wen Ye, and Yan Liu. Tempo:
591 Prompt-based generative pre-trained transformer for time series forecasting, 2024. URL [https://](https://arxiv.org/abs/2310.04948)
592 arxiv.org/abs/2310.04948. 8
- 593 Cristian Challu, Kin G. Olivares, Boris N. Oreshkin, Federico Garza, Max Mergenthaler, and Artur
Dubrawski. N-HiTS: Neural hierarchical interpolation for time series forecasting. In *Proceedings*
of the 37th AAAI Conference on Artificial Intelligence, 2022. 4, 5, 10

- 594 Ching Chang, Wen-Chih Peng, and Tien-Fu Chen. Llm4ts: Two-stage fine-tuning for time-series
595 forecasting with pre-trained llms. *arXiv preprint arXiv:2308.08469*, 2023. 8
- 596
- 597 Ben Cohen, Emaad Khwaja, Kan Wang, Charles Masson, Elise Ramé, Youssef Doubli, and Oth-
598 mane Abou-Amal. Toto: Time series optimized transformer for observability. *arXiv preprint*
599 *arXiv:2407.07874*, 2024. 3, 8
- 600 Yezhen Cong, Samar Khanna, Chenlin Meng, Patrick Liu, Erik Rozi, Yutong He, Marshall
601 Burke, David B. Lobell, and Stefano Ermon. SatMAE: Pre-training transformers for tempo-
602 ral and multi-spectral satellite imagery. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave,
603 and Kyunghyun Cho (eds.), *Advances in Neural Information Processing Systems*, 2022. URL
604 <https://openreview.net/forum?id=WBhqzPF6KYH>. 1, 2, 3, 7, 17
- 605 Hugo Dalla-Torre, Liam Gonzalez, Javier Mendoza Revilla, Nicolas Lopez Carranza, Adam Henryk
606 Grzywaczewski, Francesco Oteri, Christian Dallago, Evan Trop, Hassan Sirelkhatim, Guillaume
607 Richard, Marcin Skwark, Karim Beguir, Marie Lopez, and Thomas Pierrot. The nucleotide trans-
608 former: Building and evaluating robust foundation models for human genomics. *bioRxiv*, 2023.
609 doi: 10.1101/2023.01.11.523679. URL [https://www.biorxiv.org/content/early/
610 2023/01/15/2023.01.11.523679](https://www.biorxiv.org/content/early/2023/01/15/2023.01.11.523679). 1, 3, 6, 15, 17
- 611
- 612 Abhimanyu Das, Weihao Kong, Rajat Sen, and Yichen Zhou. A decoder-only foundation model for
613 time-series forecasting. *arXiv preprint arXiv:2310.10688*, 2023. 3
- 614 Abhimanyu Das, Weihao Kong, Rajat Sen, and Yichen Zhou. A decoder-only foundation model for
615 time-series forecasting, 2024. URL <https://arxiv.org/abs/2310.10688>. 8
- 616
- 617 Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep
618 bidirectional transformers for language understanding. In *North American Chapter of the Associ-
619 ation for Computational Linguistics*, 2019. URL [https://api.semanticscholar.org/
620 CorpusID:52967399](https://api.semanticscholar.org/CorpusID:52967399). 1, 2
- 621 Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas
622 Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszko-
623 reit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at
624 scale, 2021. URL <https://arxiv.org/abs/2010.11929>. 1
- 625 Anthony Fuller, Koreen Millard, and James R Green. Croma: Remote sensing representations
626 with contrastive radar-optical masked autoencoders. In *Thirty-seventh Conference on Neural
627 Information Processing Systems*, 2023. 1, 7, 16, 17
- 628
- 629 Mononito Goswami, Konrad Szafer, Arjun Choudhry, Yifu Cai, Shuo Li, and Artur Dubrawski.
630 Moment: A family of open time-series foundation models, 2024. URL [https://arxiv.
631 org/abs/2402.03885](https://arxiv.org/abs/2402.03885). 2, 8
- 632 Xin Guo, Jiangwei Lao, Bo Dang, Yingying Zhang, Lei Yu, Lixiang Ru, Liheng Zhong, Ziyuan
633 Huang, Kang Wu, Dingxiang Hu, et al. Skysense: A multi-modal remote sensing foundation
634 model towards universal interpretation for earth observation imagery. In *Proceedings of the
635 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 27672–27683, 2024.
636 7
- 637
- 638 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recog-
639 nition, 2015. URL <https://arxiv.org/abs/1512.03385>. 3
- 640 Patrick Helber, Benjamin Bischke, Andreas Dengel, and Damian Borth. Eurosat: A novel dataset
641 and deep learning benchmark for land use and land cover classification, 2019. URL [https://
642 //arxiv.org/abs/1709.00029](https://arxiv.org/abs/1709.00029). 7, 15
- 643 R.J. Hyndman and Y. Khandakar. Automatic time series forecasting: The forecast package for R.
644 *Journal of Statistical Software*, 27(3), 2008. 4, 5, 8, 10
- 645
- 646 Amanda A. Boatswain Jacques, Abdoulaye Baniré Diallo, and Etienne Lord. The canadian cropland
647 dataset: A new land cover dataset for multitemporal deep learning classification in agriculture,
2023. URL <https://arxiv.org/abs/2306.00114>. 7, 15

- 648 Johannes Jakubik, Sujit Roy, C. E. Phillips, Paolo Fraccaro, Denys Godwin, Bianca Zadrozny,
649 Daniela Szwarcman, Carlos Gomes, Gabby Nyirjesy, Blair Edwards, Daiki Kimura, Naomi Si-
650 mumba, Linsong Chu, S. Karthik Mukkavilli, Devyani Lambhate, Kamal Das, Ranjini Bangalore,
651 Dário A. B. Oliveira, Michal Muszynski, Kumar Ankur, Muthukumaran Ramasubramanian, Ik-
652 sha Gurung, Sam Khallaghi, Hanxi Li, Michael Cecil, Maryam Ahmadi, Fatemeh Kordi, Hamed
653 Alemohammad, Manil Maskey, Raghu K. Ganti, Kommy Weldemariam, and Rahul Ramachan-
654 dran. Foundation models for generalist geospatial artificial intelligence. *CoRR*, abs/2310.18660,
655 2023. 7
- 656 Yanrong Ji, Zhihan Zhou, Han Liu, and Ramana V Davuluri. DNABERT: pre-trained Bidirectional
657 Encoder Representations from Transformers model for DNA-language in genome. *Bioinformat-
658 ics*, 37(15):2112–2120, 02 2021. ISSN 1367-4803. doi: 10.1093/bioinformatics/btab083. URL
659 <https://doi.org/10.1093/bioinformatics/btab083>. 1, 2, 3, 6
- 660 Ming Jin, Shiyu Wang, Lintao Ma, Zhixuan Chu, James Y. Zhang, Xiaoming Shi, Pin-Yu Chen,
661 Yuxuan Liang, Yuan-Fang Li, Shirui Pan, and Qingsong Wen. Time-llm: Time series forecasting
662 by reprogramming large language models. In *The Twelfth International Conference on Learning
663 Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenReview.net, 2024. URL
664 <https://openreview.net/forum?id=Unb5CVPtae>. 8
- 665 Denis Kwiatkowski, Peter C.B. Phillips, Peter Schmidt, and Yongcheol Shin. Testing the null hy-
666 pothesis of stationarity against the alternative of a unit root: How sure are we that economic time
667 series have a unit root? *Journal of Econometrics*, 54(1):159–178, 1992. 5
- 668 Alexandre Lacoste, Nils Lehmann, Pau Rodriguez, Evan Sherwin, Hannah Kerner, Björn Lütjens,
669 Jeremy Irvin, David Dao, Hamed Alemohammad, Alexandre Drouin, et al. Geo-bench: Toward
670 foundation models for earth monitoring. *Advances in Neural Information Processing Systems*, 36,
671 2024. 3, 7, 15
- 672 Colin Lea, Rene Vidal, Austin Reiter, and Gregory D Hager. Temporal convolutional networks:
673 A unified approach to action segmentation. In *Computer Vision—ECCV 2016 Workshops: Am-
674 sterdam, The Netherlands, October 8-10 and 15-16, 2016, Proceedings, Part III 14*, pp. 47–54.
675 Springer, 2016. 3, 4
- 676 Liam Li, Kevin Jamieson, Afshin Rostamizadeh, Ekaterina Gonina, Jonathan Ben-Tzur, Moritz
677 Hardt, Benjamin Recht, and Ameet Talwalkar. A system for massively parallel hyperparameter
678 tuning. *Proceedings of Machine Learning and Systems*, 2:230–246, 2020. 4, 5
- 679 Percy Liang, Rishi Bommasani, Tony Lee, Dimitris Tsipras, Dilara Soylu, Michihiro Yasunaga, Yian
680 Zhang, Deepak Narayanan, Yuhuai Wu, Ananya Kumar, et al. Holistic evaluation of language
681 models. *arXiv preprint arXiv:2211.09110*, 2022. 1
- 682 Hanxiao Liu, Karen Simonyan, and Yiming Yang. Darts: Differentiable architecture search. *arXiv
683 preprint arXiv:1806.09055*, 2018. 4
- 684 Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo.
685 Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the
686 IEEE/CVF international conference on computer vision*, pp. 10012–10022, 2021. 1, 7
- 687 Zhuang Liu, Hanzi Mao, Chao-Yuan Wu, Christoph Feichtenhofer, Trevor Darrell, and Saining Xie.
688 A convnet for the 2020s. In *Proceedings of the IEEE/CVF conference on computer vision and
689 pattern recognition*, pp. 11976–11986, 2022. 3, 4
- 690 Oscar Manas, Alexandre Lacoste, Xavier Giró-i Nieto, David Vazquez, and Pau Rodriguez. Seasonal
691 contrast: Unsupervised pre-training from uncurated remote sensing data. In *Proceedings of the
692 IEEE/CVF International Conference on Computer Vision*, pp. 9414–9423, 2021. 7
- 693 Matías Mendieta, Boran Han, Xingjian Shi, Yi Zhu, and Chen Chen. Towards geospatial foundation
694 models via continual pretraining. In *Proceedings of the IEEE/CVF International Conference on
695 Computer Vision*, pp. 16806–16816, 2023. 1, 3, 7, 17

- 702 Eric Nguyen, Michael Poli, Marjan Faizi, Armin Thomas, Michael Wornow, Callum Birch-Sykes,
703 Stefano Massaroli, Aman Patel, Clayton Rabideau, Yoshua Bengio, et al. Hyenadna: Long-range
704 genomic sequence modeling at single nucleotide resolution. *Advances in neural information*
705 *processing systems*, 36, 2024. 1, 3, 6
- 706 Alec Radford and Karthik Narasimhan. Improving language understanding by generative pre-
707 training. 2018. URL <https://api.semanticscholar.org/CorpusID:49313245>.
708 1
- 709 Colorado J Reed, Ritwik Gupta, Shufan Li, Sarah Brockman, Christopher Funk, Brian Clipp, Kurt
710 Keutzer, Salvatore Candido, Matt Uyttendaele, and Trevor Darrell. Scale-mae: A scale-aware
711 masked autoencoder for multiscale geospatial representation learning. In *Proceedings of the*
712 *IEEE/CVF International Conference on Computer Vision*, pp. 4088–4099, 2023. 7
- 713 Nicholas Roberts, Samuel Guo, Cong Xu, Ameet Talwalkar, David Lander, Lvfang Tao, Linhang
714 Cai, Shuaicheng Niu, Jianyu Heng, Hongyang Qin, Minwen Deng, Johannes Hog, Alexander Pf-
715 efferle, Sushil Ammanaghatta Shivakumar, Arjun Krishnakumar, Yubo Wang, Rhea Sukthanker,
716 Frank Hutter, Euxhen Hasanaj, Tien-Dung Le, Mikhail Khodak, Yuriy Nevmyvaka, Kashif Rasul,
717 Frederic Sala, Anderson Schneider, Junhong Shen, and Evan Randall Sparks. AutoML Decathlon:
718 diverse tasks, modern methods, and efficiency at scale. In *Advances in Neural Information Pro-*
719 *cessing Systems: Competition Track*, 2021a. 4
- 720 Nicholas Roberts, Mikhail Khodak, Tri Dao, Liam Li, Christopher Ré, and Ameet Talwalkar. Re-
721 thinking neural operations for diverse tasks. *Advances in Neural Information Processing Systems*,
722 34:15855–15869, 2021b. 4
- 723 Esther Rolf, Konstantin Klemmer, Caleb Robinson, and Hannah Kerner. Mission critical–satellite
724 data is a distinct modality in machine learning. *arXiv preprint arXiv:2402.01444*, 2024. 3, 7
- 725 Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomed-
726 ical image segmentation, 2015. URL <https://arxiv.org/abs/1505.04597>. 4
- 727 Junhong Shen, Mikhail Khodak, and Ameet Talwalkar. Efficient architecture search for diverse
728 tasks. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2022. 4, 5
- 729 Junhong Shen, Liam Li, Lucio M Dery, Corey Staten, Mikhail Khodak, Graham Neubig, and Ameet
730 Talwalkar. Cross-modal fine-tuning: Align then refine. In *International Conference on Machine*
731 *Learning*, pp. 31030–31056. PMLR, 2023. 4
- 732 Gencer Sumbul, Marcela Charfuelan, Begum Demir, and Volker Markl. Bigearthnet: A large-
733 scale benchmark archive for remote sensing image understanding. In *IGARSS 2019 - 2019 IEEE*
734 *International Geoscience and Remote Sensing Symposium*. IEEE, July 2019. doi: 10.1109/igarss.
735 2019.8900532. URL <http://dx.doi.org/10.1109/IGARSS.2019.8900532>. 7, 15
- 736 Jingmin Sun, Yuxuan Liu, Zecheng Zhang, and Hayden Schaeffer. Towards a foundation model
737 for partial differential equation: Multi-operator learning and extrapolation. *arXiv preprint*
738 *arXiv:2404.12355*, 2024. 2
- 739 Renbo Tu, Nicholas Roberts, Mikhail Khodak, Junhong Shen, Frederic Sala, and Ameet Talwalkar.
740 NAS-Bench-360: Benchmarking diverse tasks for neural architecture search. In *Advances in*
741 *Neural Information Processing Systems: Datasets and Benchmarks Track*, 2022. 4
- 742 Yuxuan Wang, Haixu Wu, Jiayang Dong, Yong Liu, Mingsheng Long, and Jianmin Wang. Deep
743 time series models: A comprehensive survey and benchmark. 2024. 16
- 744 Gerald Woo, Chenghao Liu, Akshat Kumar, Caiming Xiong, Silvio Savarese, and Doyen Sa-
745 hoo. Unified training of universal time series forecasting transformers, 2024. URL <https://arxiv.org/abs/2402.02592>. 8
- 746 Sergey Zagoruyko and Nikos Komodakis. Wide residual networks, 2017. URL <https://arxiv.org/abs/1605.07146>. 4

- 756 Ailing Zeng, Muxi Chen, Lei Zhang, and Qiang Xu. Are transformers effective for time series
757 forecasting? In *Proceedings of the AAAI conference on artificial intelligence*, volume 37, pp.
758 11121–11128, 2023. 3, 8
- 759 Ruijie Zhao, Mingwei Zhan, Xianwen Deng, Yanhao Wang, Yijun Wang, Guan Gui, and Zhi Xue.
760 Yet another traffic classifier: A masked autoencoder based traffic transformer with multi-level
761 flow representation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37,
762 pp. 5420–5427, 2023. 2
- 764 Jian Zhou and Olga G Troyanskaya. Predicting effects of noncoding variants with deep learn-
765 ing-based sequence model. *Nature Methods*, 12(10):931–934, Aug 2015. doi: 10.1038/nmeth.
766 3547. 6
- 767 Tian Zhou, PeiSong Niu, Xue Wang, Liang Sun, and Rong Jin. One fits all:power general time series
768 analysis by pretrained lm, 2023a. URL <https://arxiv.org/abs/2302.11939>. 8
- 770 Zhihan Zhou, Yanrong Ji, Weijian Li, Pratik Dutta, Ramana Davuluri, and Han Liu. Dnabert-2:
771 Efficient foundation model and benchmark for multi-species genome, 2023b. 1, 3, 6

772 APPENDIX

773 A FURTHER EXPERIMENTAL DETAILS

774 A.1 HYPERPARAMETER SEARCH SPACE

775 Table 4: Hyperparameter Search Space

Hyperparameter	Search Space	Type of Search Space
random_seed	[0, 500]	Integer
lr	$[10^{-5}, 5 \times 10^{-1}]$	Log Uniform
drop_rate	{0, 0.05, 0.1}	Discrete
weight_decay	$[5 \times 10^{-7}, 5 \times 10^{-3}]$	Log Uniform
momentum	[0.9, 1]	Uniform

788 A.2 DATASETS

789 **Genomics** For the Genomics domain, we use the 18 classification tasks from the Nucleotide Trans-
790 former benchmark Dalla-Torre et al. (2023) that has widely been used for other genomics FMs. The
791 benchmark datasets consist of nucleotide base sequences ranging from 200 to 600 bases in length.
792 It provides a realistic and biological meaningful benchmark across four main categories: promotor
793 (human/mouse), enhancer (human), splice site (SS) (human/multispecies) and histone modification
794 (yeast). Within the benchmark, the `enhancers_types` and `splice_sites_all` datasets are
795 classification tasks with three classes each, while the remaining datasets are binary classification
796 tasks.

797 **Satellite image** In the satellite imaging domain, we aim to conduct evaluations with real-world rele-
798 vance to Earth science. To achieve this, we include a variety of data from different sources to cover
799 a diverse range of tasks, such as brick kiln identification, deforestation prediction, and photovoltaic
800 monitoring. We utilize five classification tasks provided by the GeoBench dataset (Lacoste et al.,
801 2024), a recently developed benchmark that offers a clean and carefully curated collection of tasks
802 specifically designed for satellite imaging. In addition to GeoBench, we evaluate our model on three
803 additional datasets (Helber et al., 2019; Jacques et al., 2023; Sumbul et al., 2019) commonly used in
804 the literature as benchmarks for this domain. This brings the total to eight datasets, encompassing a
805 wide range of features. These tasks vary in complexity, with single-class classification ranging from
806 binary to 62-class problems, as well as two multilabel classification tasks. The datasets are further
807 characterized by diverse input channels, ranging from 3 RGB channels to 18 channels that integrate
808 data from both Sentinel-1 and Sentinel-2 formats.

810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829

Dataset	# of classes	# of samples train	test	Maximum sequence length
enhancers	2	14968	400	200
enhancers_types	3	14968	400	200
promoter_all	2	53276	5920	300
promoter_no_tata	2	47767	5299	300
promoter_tata	2	5509	621	300
splice_sites_acceptors	2	19961	2218	600
splice_sites_all	3	27000	3000	400
splice_sites_donors	2	19775	2198	600
H3	2	13468	1497	500
H3K14ac	2	29743	3305	500
H3K36me3	2	31392	3488	500
H3K4me1	2	28509	3168	500
H3K4me2	2	27614	3069	500
H3K4me3	2	25953	2884	500
H3K79me3	2	25953	2884	500
H3K9ac	2	25003	2779	500
H4	2	13140	1461	500
H4ac	2	30685	3410	500

Table 5: Statistics for Genomics datasets

830
831
832
833

For Geo-Bench datasets, we do not use any `mixup` and `cutmix` augmentations. For other datasets, we universally use `mixup` = 0.8, `cutmix` = 1.0, and a switch probability of 0.5. Following (Fuller et al., 2023), we use only 10% of training set from BigEarthNet and fMoW-Sentinel while using the full evaluation set for validation.

834
835
836
837
838
839
840
841
842
843
844
845
846
847
848

Dataset	Image Size	# of classes	# of samples train	val	test	# of channels
m-bigearthnet	120 × 120	43	20000	1000	1000	12
m-brickkiln	64 × 64	2	15063	999	999	13
m-so2sat	32 × 32	17	19992	986	986	18
m-forestnet	332 × 332	12	6464	989	993	6
m-pv4ger	320 × 320	2	11814	999	999	3
BigEarthNet	120 × 120	19	31166	103944	103728	12
EuroSAT	64 × 64	13	16200	10800	5400	13
Canadian Cropland	120 × 120	10	53884	11414	11674	12
fMoW-Sentinel	96 × 96	62	71287	84939	84966	13

Table 6: Statistics for Satellite datasets

849
850
851
852
853
854
855
856
857

Time series In the time series domain, we focus on the long horizon forecasting task. We use a subset of the common benchmark datasets for evaluating models across different domains (ETT, Electricity, Weather, Illness, Traffic, Exchange Rate) (Wang et al., 2024), specifically, the ETT, Waether, Electricity, Illness, and Traffic datasets. Note that the ETT dataset is actually a collection of four series: ETTh1, ETTh2, ETTm1, and ETTm2; we follow the rest of the literature in treating each series as a separate dataset. Each dataset contains measurements of one or more channels at evenly spaced time steps.

858
859

B RESULTS AND IMPLEMENTATION DETAILS

860
861
862
863

ASHA tuning details Following the architecture search, we perform hyperparameter tuning using ASHA. The hyperparameter search space includes learning rate, weight decay, momentum, drop rate, and random seed for model initialization. We define a continuous search space, with further specific details provided in Appendix A.1. Using ASHA, we evaluate 200 sample configurations

864
865
866
867
868
869
870
871
872
873
874

Dataset	# of channels	# of samples		
		train	val	test
ETTh1	7	8033	2785	2785
ETTh2	7	8033	2785	2785
ETTm1	7	33953	11425	11425
ETTm2	7	33953	11425	11425
Weather	21	36280	5175	10444
Electricity	321	17805	2537	5165
ILI	7	69	2	98
Traffic	862	11673	1661	3413

875
876
877

Table 7: Statistics for Time Series datasets

878
879

over a maximum of 20 epochs, using a reduction factor of 2. The low-performing configurations are pruned based on their validation scores.

880
881
882
883
884

Before retraining the final model, we load the model checkpoint corresponding to the optimal hyperparameter configuration. The model is then trained for 200 epochs on the training data, with the best-performing checkpoint selected based on validation performance. This process is repeated for each backbone architecture, and the best-performing backbone is selected using the validation score. Finally, the checkpoint for the selected backbone is evaluated on the test set to obtain the final score.

885
886
887
888
889
890
891
892
893

Genomics As the score presented are taken from [Dalla-Torre et al. \(2023, Supplementary Table 6\)](#), we do realize that in their reported table, all the promoter and splice sites tasks are mislabeled with some histone modification tasks. In order to get the results, we infer an order for the mislabeled datasets with most confidence. We also include all FMs listed on the leaderboard in our evaluation for comprehensive comparison. In alignment with the leaderboard, we apply a 0.1 validation split for DASHA during our evaluation. Additionally, we use an architecture set that includes both Wide ResNet and UNet for the search with DASHA on these datasets. We use batch size= 128 for all datasets, and cross entropy loss for all the training and finetuning. Individual scores for each task in the benchmark are provided in Tables 9 and 8.

894
895
896
897
898
899
900
901
902
903
904
905

Model	Regulatory Elements					RNA Production		
	enhancers	enhancers types	promoter all	promoter no_tata	promoter tata	splice_sites acceptors	splice_sites all	splice_sites donors
NT-Multispecies-v2 (500M)	0.559	0.438	0.976	0.976	0.965	0.981	0.984	0.987
NT-Multispecies (2.5B)	0.545	0.444	0.975	0.977	0.959	0.986	0.982	0.987
NT-1000G (500M)	0.509	0.395	0.951	0.951	0.936	0.965	0.968	0.971
NT-1000G (2.5B)	0.546	0.432	0.965	0.967	0.957	0.98	0.976	0.979
HyenaDNA-32K	0.489	0.352	0.956	0.954	0.939	0.96	0.962	0.957
HyenaDNA-1K	0.52	0.403	0.959	0.959	0.944	0.959	0.956	0.947
Enformer	0.454	0.312	0.955	0.955	0.959	0.915	0.847	0.906
DNABERT-2	0.525	0.423	0.972	0.972	0.955	0.975	0.939	0.963
DNABERT-1	0.495	0.367	0.961	0.962	0.956	-	0.975	-
Wide ResNet	0.525	0.416	0.915	0.914	0.890	0.659	0.279	0.608
UNet	0.490	0.366	0.910	0.919	0.896	0.944	0.265	0.618
DASHA	0.527	0.432	0.923	0.920	0.903	0.959	0.972	0.961

906
907

Table 8: Regulatory Elements and RNA Production Downstream Tasks

908
909
910
911
912
913
914
915

Satellite Imaging Training on satellite datasets requires relatively large computational resources due to the high number of channels and the size of the datasets. To ensure a fair comparison, we fine-tuned all the foundation models ourselves by sweeping across a fixed set of base learning rates $[5e - 3, 2e - 3, 2e - 4, 4e - 5]$. We then calculate the actual learning rate from base learning rate following previous work by $lr = base_lr \cdot \frac{batch_size}{256}$. This approach ensures that approximately the same amount of resources were used as during the DASHA tuning process, allowing for a balanced evaluation of model performance.

916
917

We closely followed the reported evaluation processes from previous studies on FMs ([Cong et al., 2022](#); [Fuller et al., 2023](#); [Mendieta et al., 2023](#)). These models do not employ a validation set for hyperparameter tuning or model selection, and we adhered to this same approach when fine-tuning

Model	H3	H3K14ac	H3K36me3	H3K4me1	H3K4me2	H3K4me3	H3K79me3	H3K9ac	H4	H4ac
NT-Multispecies-v2 (500M)	0.786	0.549	0.624	0.55	0.32	0.406	0.63	0.567	0.799	0.496
NT-Multispecies (2.5B)	0.793	0.538	0.618	0.541	0.324	0.408	0.623	0.547	0.808	0.492
NT-1000G (500M)	0.736	0.381	0.468	0.38	0.26	0.235	0.562	0.479	0.755	0.342
NT-1000G (2.5B)	0.754	0.453	0.53	0.418	0.278	0.311	0.574	0.491	0.787	0.408
HyenaDNA-32K	0.747	0.405	0.479	0.387	0.276	0.291	0.567	0.472	0.761	0.379
HyenaDNA-1K	0.781	0.608	0.614	0.512	0.455	0.55	0.669	0.586	0.763	0.564
Enformer	0.724	0.284	0.345	0.291	0.207	0.156	0.498	0.415	0.735	0.275
DNABERT-2	0.785	0.515	0.591	0.512	0.333	0.353	0.615	0.545	0.797	0.465
DNABERT-1	0.763	0.403	0.474	0.396	0.282	0.258	0.578	0.505	0.784	0.359
Wide ResNet	0.798	0.667	0.670	0.554	0.541	0.660	0.706	0.620	0.754	0.657
UNet	0.797	0.647	0.482	0.541	0.553	0.292	0.562	0.624	0.760	0.389
DASHA	0.790	0.683	0.630	0.528	0.640	0.714	0.721	0.709	0.776	0.744

Table 9: Histone Modification Downstream Tasks

Model	Average Score \uparrow	Average Rank \downarrow	Mean %Imp. \uparrow	Median %Imp. \uparrow
Foundation Model				
CROMA-base	77.98	3.57	6.95	5.30
CROMA-large	79.06	2.14	8.84	6.26
SatMAE-base	77.71	6.00	6.74	4.56
SatMAE-large	78.44	4.07	7.87	6.75
GFM	76.95	5.57	5.40	4.08
SwinT-Base	76.12	6.50	3.98	1.43
Supervised Methods				
ResNet50	73.25	9.00	-0.27	-0.38
Wide ResNet	73.90	8.00	0.00	0.00
UNet	75.29	6.57	2.33	0.87
DASHA	78.00	3.57	7.02	5.16

Figure 5: Aggregated metrics on Satellite Imaging datasets, excluding the m-pv4ger and m-forestnet datasets, where input channels do not match the model requirements for FMs

the FMs. However, for DASHA, since we performed extensive hyperparameter optimization over a large search space, we used a validation set to ensure fair and accurate comparisons between DASHA and the FMs. This is a less favorable setting for DASHA, as it relies on extensive hyperparameter tuning, but we demonstrate that, even under these conditions, DASHA matches the performance of the FMs.

It is also important to note that SatMAE only accepts 3-channel and 12-channel inputs, while CROMA is limited to 12-channel inputs. GeoBench, however, includes a wide range of tasks with varying numbers of input channels, ranging from 3 to 18. Despite these differences, we include all datasets in our evaluation because they are valuable benchmarks in the satellite image domain, and it is crucial for FMs in this field to generalize across diverse datasets. For datasets where the input size does not match the model requirements, we pad missing channels with zeros and prune any extra channels. However, to ensure a fair comparison, in addition to reporting the average scores across all datasets, we also provide average scores excluding m-pv4ger and m-forestnet, where missing channels may affect the performance of the FMs. The performance profile and aggregate scores excluding m-pv4ger and m-forestnet are presented in Figure 5.

For training and finetuning, we universally use `batch_size = 16` and loss function as cross entropy with 0.1 label smoothing for single label classification, and multi-label soft margin loss for multilabel classification. Individual scores for each task are provided in Table 10.

Time series The long horizon forecasting task for a time series can be summarized as follows: at every timestep t , take the L historical observations at times $(t - L + 1, \dots, t)$ for each channel and predict the next H observations $(t + 1, \dots, t + H)$ for every channel. Following the literature, we evaluate models on $H \in \{24, 36, 48, 60\}$ for the Illness dataset and $H \in \{96, 192, 336, 720\}$. For most methods, we report results when $L = 512$.

In addition to the results for DASHA, Auto-AR, DLinear, and FMs, we evaluate the performance of two simple baselines

Model	Average	m-bigearthnet	m-brickkiln	m-so2sat	m-forestnet	m-pv4ger	BigEarth Net	EuroSAT	Canadian Cropland	fMoW Sentinel
CROMA-base	77.39	72.07	98.99	60.04	54.07	96.6	86.94	98.81	75.87	53.14
CROMA-large	78.03	73.36	99.01	59.22	51.96	96.9	87.98	98.98	76.56	58.32
SatMAE-base	76.99	72.3	98.22	54.56	51.89	97.0	86.04	98.69	74.64	59.55
SatMAE-large	77.75	73.82	98.6	55.79	53.7	96.92	86.75	98.86	75.38	59.89
GFM	77.18	71.97	98.35	57.52	59.38	96.54	85.93	99.02	72.03	53.84
ResNet50	73.76	60.31	98.4	51.83	54.29	96.79	79.16	98.23	72	52.84
SwinT-base	76.69	70.14	98.81	56.49	59.78	97.54	85.91	98.99	70.15	52.37
Wide ResNet	73.97	69.15	98.95	49.04	52.14	96.34	80.48	98.6	72.05	49.00
UNet	75.73	69.89	98.6	56.87	57.18	97.39	83.9	98.99	72.68	46.11
DASHA	77.85	72.72	98.92	56.28	57.24	97.4	86.09	99.07	75.69	57.20

Table 10: Satellite Imaging Tasks

Model	ETTth1				ETTth2				ETTh1				ETTh2			
	96	192	336	720	96	192	336	720	96	192	336	720	96	192	336	720
Horizon	96	192	336	720	96	192	336	720	96	192	336	720	96	192	336	720
GPT4TS (OFA)	0.376	0.416	0.442	0.477	0.285	0.354	0.373	0.406	0.292	0.332	0.366	0.417	0.173	0.229	0.286	0.378
LLM4TS	0.371	0.403	0.42	0.422	0.269	0.328	0.353	0.383	0.285	0.324	0.353	0.408	0.165	0.22	0.268	0.35
MOMENT	0.387	0.41	0.422	0.454	0.288	0.349	0.369	0.403	0.293	0.326	0.352	0.405	0.17	0.227	0.275	0.363
TEMPO (Zero Shot)	0.4	0.426	0.441	0.443	0.301	0.355	0.379	0.409	0.438	0.461	0.515	0.591	0.185	0.243	0.309	0.386
TimesFM (Zero Shot)	0.421	0.472	0.51	0.514	0.326	0.399	0.434	0.451	0.357	0.411	0.441	0.507	0.205	0.294	0.367	0.473
Moirai (Zero Shot)	0.375	0.399	0.412	0.413	0.281	0.34	0.362	0.38	0.404	0.435	0.462	0.49	0.205	0.261	0.319	0.415
Toto (Zero Shot)	0.307	0.329	0.396	0.419	0.093	0.135	0.16	0.294	0.306	0.328	0.39	0.463	0.2	0.269	0.264	0.354
ARIMA	0.646	0.704	0.732	0.738	0.324	0.411	0.456	0.462	1.131	1.172	1.197	1.231	0.225	0.298	0.37	0.478
AR (d=0)	0.358	0.39	0.41	0.424	0.271	0.334	0.361	0.395	0.299	0.336	0.368	0.426	0.163	0.218	0.271	0.366
DLinear	0.375	0.405	0.439	0.472	0.289	0.383	0.448	0.605	0.299	0.335	0.369	0.425	0.167	0.224	0.281	0.397
Auto-AR	0.357	0.39	0.41	0.422	0.269	0.332	0.359	0.394	0.299	0.336	0.368	0.426	0.163	0.218	0.271	0.367
DASHA	0.369	0.401	0.430	0.478	0.284	0.377	0.396	0.745	0.305	0.335	0.367	0.418	0.169	0.224	0.290	0.378
Model	Weather				Electricity				ILI				Traffic			
	96	192	336	720	96	192	336	720	96	192	336	720	96	192	336	720
Horizon	96	192	336	720	96	192	336	720	96	192	336	720	96	192	336	720
GPT4TS (OFA)	0.162	0.204	0.254	0.326	0.139	0.153	0.169	0.206	2.063	1.868	1.79	1.979	0.388	0.407	0.412	0.45
LLM4TS	0.147	0.191	0.241	0.313	0.128	0.146	0.163	0.2	-	-	-	-	0.372	0.391	0.405	0.437
MOMENT	0.154	0.197	0.246	0.315	0.136	0.152	0.167	0.205	2.728	2.669	2.728	2.883	0.391	0.404	0.414	0.45
TEMPO (Zero Shot)	0.211	0.254	0.292	0.37	0.178	0.198	0.209	0.279	3.0	2.956	2.651	2.701	0.476	0.496	0.503	0.538
TimesFM (Zero Shot)	0.122	0.169	0.242	0.391	0.119	0.137	0.158	0.206	2.595	2.984	3.34	3.227	0.327	0.353	0.378	0.42
Moirai (Zero Shot)	0.173	0.216	0.26	0.32	0.205	0.22	0.236	0.27	-	-	-	-	-	-	-	-
Toto (Zero Shot)	0.18	0.235	0.252	0.356	0.124	0.138	0.155	0.211	-	-	-	-	-	-	-	-
ARIMA	0.217	0.263	0.33	0.425	1.22	1.264	1.311	1.364	5.554	6.94	7.192	6.648	1.997	2.044	2.096	2.138
AR (d=0)	0.171	0.215	0.263	0.332	0.138	0.153	0.17	0.212	2.084	2.04	2.004	2.011	0.398	0.413	0.426	0.464
DLinear	0.176	0.22	0.265	0.323	0.14	0.153	0.169	0.203	2.215	1.963	2.13	2.368	0.41	0.423	0.436	0.466
Auto-AR	0.172	0.215	0.263	0.332	0.138	0.153	0.17	0.212	2.084	2.04	2.004	2.011	0.398	0.413	0.426	0.464
DASHA	0.163	0.205	0.251	0.314	0.136	0.151	0.165	0.200	-	-	-	-	-	-	-	-

Table 11: Time Series Forecasting Tasks

1. Vanilla Autoregressive Model (Box & Jenkins, 1976): This model predicts the (scalar) value of a time series at $t + 1$ as a linear combination of the last L timesteps and a constant, i.e. $\hat{x}_{t+1} = \alpha_0 + \alpha_1 x_t + \alpha_2 x_{t-1} + \dots + \alpha_L x_{t-L+1}$ for learnable parameters $\alpha_0, \dots, \alpha_L$. We fit these parameters using standard maximum likelihood techniques.
2. ARIMA is a statistical method used for time series forecasting that combines three components: AutoRegressive (AR), Integrated (I), and Moving Average (MA). The AR component models the relationship between an observation and its lagged (past) values, assuming that past values have a linear influence on future ones. The Integrated component applies differencing to the data to remove trends or seasonality, making the time series stationary by stabilizing its mean over time. The MA component models the relationship between an observation and the residual errors from a moving average model applied to previous observations. ARIMA is characterized by three parameters: p (the number of lag observations), d (the number of differencing steps to achieve stationarity), and q (the number of lagged forecast errors). This model is particularly effective for univariate time series forecasting where patterns like trends or seasonality are present.

All results are reported on a 70/10/20 train/validation/test split for each datasets, except for the ETT datasets which have predefined splits. MSE is reported after all datasets have been scaled by the mean and variance of the training data. Both autoregressive models have only one tunable hyperparameter (number of lags). Similarly, the linear model has only one tunable hyperparameter (number of training epochs).

1026 The baselines as described can handle only univariate time series, while all of the benchmark datasets
1027 are multivariate (multiple channels). These baselines, as well as DASHA, are trained under channel
1028 independence: each channel of a time series is treated independently. While channel independence
1029 fails to take into account cross-channel dependencies, we note that developing methods that leverage
1030 cross-channel dependencies for a variable number of channels remains an open problem.

1031
1032
1033
1034
1035
1036
1037
1038
1039
1040
1041
1042
1043
1044
1045
1046
1047
1048
1049
1050
1051
1052
1053
1054
1055
1056
1057
1058
1059
1060
1061
1062
1063
1064
1065
1066
1067
1068
1069
1070
1071
1072
1073
1074
1075
1076
1077
1078
1079