
EEVEE and GATE: Finding the right benchmarks and how to run them seamlessly

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Model evaluation is a cornerstone of machine learning, guiding model design and
2 progress measurement. Designing generalizable evaluation processes remains a
3 challenge, however, partly due to the vast number of possible domain, task and
4 modality combinations and lack of knowledge of how informative they are. In
5 this paper, we propose *EEVEE* (Efficient Evaluation process Evolution Engine)¹, a
6 method that frames evaluation process design as a learning problem. By analyzing
7 a large number of evaluation metrics from diverse benchmarks and models, EEVEE
8 identifies a smaller subset of tasks with high predictive power over the full set of
9 evaluation metrics, reducing evaluation time. To find the optimal subset maximiz-
10 ing signal while minimizing GPU hours, EEVEE evaluates pre-trained models of
11 various architectures, pretraining schemes, and modalities on diverse downstream
12 tasks and datasets including image classification, segmentation, relational reason-
13 ing, zero-shot image-to-text tasks, medical classification and segmentation, video
14 classification, and regression. Our results identify three subsets of benchmarks,
15 with 8, 15 and 21 tasks, providing high quality signal for model generalization.
16 Key benchmarks selected include iWildCam, CLEVR-Math, ACDC, WinoGround,
17 CIFAR100, Fungi, and ADE20K. We structure the subsets into three tiers for
18 12, 24, and 36 GPU-hour budgets and package them into a unified, efficient, and
19 user-friendly Python framework that we built with the researcher in mind – which
20 we refer to as the GATE engine. Our experiments reveal ConvNextV2, SigLIP
21 and CLIP as top-performing model encoders, with EfficientNetV2 and ResNext50
22 excelling in medical tasks and challenging image classification, in particular in
23 Happy Whale Individual classification, ConvNet based models seem to outperform
24 transformer models by a factor of 2.5x, which is surprising. The top performing en-
25 coder being ConvNextV2 followed by CLIP seems to agree with other recent large
26 scale evaluations. We also demonstrate the framework’s versatility in fine-tuning
27 models from text and audio modalities, paving the way for future cross-modal
28 evaluations.

29 1 Introduction

30 **Increasing Complexities of Benchmarking:** As we create benchmarks for expanding model capa-
31 bility evaluation, the growing number and complexity of these benchmarks inadvertently complicates
32 evaluation, requiring more resources like engineering, computation, and research time. Consequently,
33 prioritizing which benchmarks to use becomes challenging. The high costs and longer wait times of
34 newer, complex benchmarks often deter their adoption, leading researchers to rely on older, simpler
35 benchmarks. This risks missing valuable insights from innovative ideas that may underperform on

¹Pronounced as /'i:vi:/ EE-vee

36 simpler benchmarks but have broader applicability, while promoting incremental improvements that
37 overfit to simpler benchmarks but underperform in comprehensive evaluations.

38 To illustrate the mounting increase in available benchmarks, we can look at the historical benchmarks
39 in deep learning. Few benchmarks have had as much impact as ImageNet [29], which remains a
40 rich resource for model training and evaluation, particularly in visuo-linguistic models. As key
41 capabilities for deep neural networks were discovered, more benchmarks were generated to measure
42 and stimulate progress in those areas. In natural language processing, the GLUE benchmark [65],
43 SQuAD [45], and CoNLL-2003 [48] have been instrumental. In audio processing, LibriSpeech [39],
44 TIMIT [15], and VCTK [68] are widely used. For machine translation, WMT [3], IWSLT [22], and
45 Europarl [25] have driven advancements. Relational reasoning has been advanced by benchmarks
46 such as CLEVR [23], bAbI [66], and RAVEN [71]. In segmentation, PASCAL VOC [14], Cityscapes
47 [8], and COCO [33] remain crucial. Large language models are often evaluated using benchmarks
48 like SuperGLUE [64], LAMBADA [40], and MMLU [19]. Vision-language models are typically
49 evaluated using benchmarks such as VQA [1], Visual7W [76], and Flickr30k [42].

50 As a result, a researcher has to choose from all these options, and even more, and then find a
51 way to unify and experiment with their models across all of them. The lack of unification, and
52 the lack of guarantees for their generalization signal, quickly becomes a kind of “evaluation hell”,
53 where researchers waste a lot of time just doing redundant things like fixing the same bugs to
54 download datasets, preprocess them etc, while at the same time not having any real signal as to which
55 benchmarks are more informative, other than just knowing what has been used the most – which is
56 usually a function of popularity, and not real informativeness. To elaborate, the adoption of complex
57 evaluation processes that could enhance research efficiency and impact is often hindered by the
58 engineering effort required to evaluate machine learning models. Researchers must create involved
59 pipelines across multiple datasets demanding high data engineering efforts, develop task-specific
60 adapters, and derive nuanced training recipes, which is time-consuming. As a result, researchers
61 often revert to simpler evaluation strategies instead of comprehensive assessments.

62 A good benchmark should alleviate these burdens by automating dataset handling, integrating task
63 adapters, optimizers, schedulers, and logging mechanisms seamlessly. It should provide broad and
64 meaningful signals with minimal GPU time, accommodating various computational budgets, ensuring
65 inclusivity. Furthermore, an increasingly important factor for a robust modern benchmark engine
66 is its support for multi-modal learning and early fusion techniques. AI systems must seamlessly
67 integrate and reason across multiple modalities, such as text, images, audio, and more. Multi-modal
68 learning enhances self-supervised learning opportunities and provides inherent supervision through
69 natural alignments, like audio-visual synchronization in videos. Early fusion, where data from
70 different modalities is combined at the initial stages of processing, allows models to leverage shared
71 representations, improving generalization and reasoning capabilities across varied tasks and domains.
72 These key desiderata are what motivates the production of this work.

73 With the desiderata in mind, we next introduce EEVEE, a methodology developed for building
74 high-signal low-cost evaluation routines, and GATE, the resulting benchmark that is designed to
75 be extensible, readable, flexible, modular and robust, supported by a new efficient, easy to use
76 framework.

77 **EEVEE, Learning Optimal Benchmarks:** The ability to find which benchmarks offer the most
78 signal with respect to a given goal, such that we can optimize our compute time, research iteration
79 speed, and engineering time is increasingly crucial. In this work, rather than just manually designing
80 a new set of benchmarks, we propose a methodology, called *EEVEE (Empirical Evaluation process
81 Evolution Engine)* that frames evaluation design as a learning problem and then leverages machine
82 learning to automate the discovery and refinement of evaluation processes.

83 More specifically, EEVEE operates by taking in a large set of performance metrics from diverse
84 models applied across various benchmarks and identifies a smaller subset of benchmarks with high
85 predictive power over the entire set. EEVEE achieves this through two main components: (a) an
86 evolutionary algorithm to optimize the selection of benchmark combinations based on a computed
87 score, and (b) a meta-model trained to predict a model’s performance on the full set of benchmarks
88 using performance metrics from a chosen subset. We parameterize the meta-model as a small
89 neural network.

90 The meta-model receives input performance metrics from a subset of benchmarks and predicts perfor-
91 mance on the full set of performance metrics. Through careful k -fold cross-validation and leveraging
92 a diverse set of models and benchmarks, EEVEE iteratively evolves benchmark combinations that
93 offer high information content with respect to the entire spectrum of benchmarks, ensuring robust,
94 efficient and comprehensive evaluation that can be targeted to computational budgets ranging from
95 more “GPU Poor” users to high-budget organizations.

96 Taking the desiderata explained above and the resulting understanding of what a good evaluation
97 engine should look like, we demonstrate the effectiveness of EEVEE by tasking it with the discovery
98 of benchmark combinations that offer good **signal-to-GPU-time** ratio, for the evaluation of **model**
99 **encoders** – also referred to as backbones, on their ability to adapt to new tasks, domains, and
100 modalities. For this purpose, we choose a pool of 20 models, varying in their pretraining schemes
101 (e.g. CLIP, DINO, ImageNet Classification), architectures (e.g. ResNets, ViTs, ConvNext) and even
102 their source modalities (e.g. Whisper, BERT), which we adapt on 31 benchmarks ranging from image
103 classification, segmentation, relational reasoning, zero-shot image-to-text tasks, medical classification
104 and segmentation, video classification, and regression, using robust fine tuning recipes, and training
105 for 10K iterations, ensuring that the signal we get is about models that are adaptable, generalizable
106 and efficient in their adaptation.

107 By applying 20 models on 31 benchmarks and employing EEVEE on their resulting metrics, we
108 identify three subsets of benchmarks, each targeted to a specific computational budget range. Some of
109 the key benchmarks that have been selected include iWildCam, CLEVR-Math, ACDC, WinoGround,
110 mini-ImageNet, Fungi, ADE20K, and dtextures. We refer to the discovered subsets as *Tiers*, and
111 assign to them identifiers for their sizes, specifically, *small* ($n=8$, 12 GPU hours), *base* ($n=15$, 24 GPU
112 hours) and *big* ($n=31$, 36 GPU hours). We package these tiers into our comprehensive benchmarking
113 suite and software framework (called *GATE*) designed for domain, task and modality transferability
114 evaluation, which facilitates the transfer of neural network encoders to different modalities, domains,
115 and tasks. *GATE*’s architecture caters to the research community, enabling straightforward replace-
116 ment of these transferable encoders with minimal effort. With these innovations, *GATE* seeks to
117 evolve the landscape of model encoder evaluation, championing a deeper understanding of transfer
118 learning and model adaptability.

119 **Contributions:** 1. We introduce *EEVEE*, a machine learning approach for selecting subsets of
120 benchmarks optimized to offer maximal predictive power over a larger benchmark set. 2. We conduct
121 a comprehensive investigation of diverse benchmarks within the space of image, image+text and
122 video modalities, pinpointing those with the highest predictive value for a model’s performance
123 in downstream tasks. We apply EEVEE to model encoder evaluation by training 20 models on 31
124 benchmarks, identifying subsets of 8, 15 and 21 benchmarks that offer high signal-to-GPU-hour ratios.
125 3. We pack the EEVEE-discovered subsets (of 8, 15 and 21 benchmarks out of 31 benchmarks) into
126 targeted benchmark packs, referred to as tiers, designed for specific compute budgets (of 12, 24 and
127 36 GPU hours) and project phases, and establish standard experimental settings for these tiers. We call
128 these collectively as the *GATE* Benchmarks. 4. We develop the *GATE* engine, a unified benchmark
129 suite and software framework that automates dataset downloading, preprocessing, and pipelining
130 for fine tuning and evaluation. *GATE* facilitates the incorporation of new model encoders, adapts
131 input modalities, fine-tunes with robust recipes, and logs critical information such as training and
132 evaluation metrics, power, energy, computational usage, task visualizations, and model gradients per
133 layer. 5. Through our extensive investigation, we identify foundation models demonstrating superior
134 transferability across diverse tasks. 6. We advocate for the inclusion of modality-shifting transfer
135 experiments in the standard evaluation process for ML researchers, supported by our experimental
136 results on the performance of existing foundation models in these benchmarks.

137 2 Related Work

138 **On the Diversity of Benchmarks:** There is a vast array of benchmark suites in machine learning.
139 To the best of our knowledge, the benchmark suites relating strongly to *GATE* are ImageNet [9],
140 VTAB [70], VLMBench [73] and WILDS [26]. ImageNet has been of tremendous importance and
141 interest to the transfer learning community. Nevertheless, there has been skepticism about overfitting
142 to such datasets resulting from implicitly qualifying models using the test set performance over
143 the years [46, 6] or the test set not being challenging enough to gauge model generalization power
144 [47]. Although ImageNet pre-training helps transfer performance to the many-shot classification
145 setting [13], it provides minimal to no gains on more challenging datasets such as fine-grained

Desiderata ↓ Benchmark →	ImageNet	VTAB	VLMBench	WILDS	GATE (Ours)
Diversity of Tasks	✓	✓	✓	✓	✓
Diversity of Domains	✓	✓	✓	✓	✓
Diversity of Modalities	✓	✓	✓	✓	✓
Automatic Dataset Download/Preparation	✓	✓	✓	✓	✓
Code allows for easy switch of encoders	✓	✓	✓	✓	✓
Optimized for fast and effective research iteration	✓	✓	✓	✓	✓
Run Time	✓	✓	✓	✓	✓
Includes Medical Domains	✓	✓	✓	✓	✓
Includes Environmental domains	✓	✓	✓	✓	✓
Tiered compute budgets	✓	✓	✓	✓	✓
GPU poor optimized	✓	✓	✓	✓	✓

Table 1: Our Desiderata (first column) VS Benchmarks (first row)

146 classification [27]. Similarly, with a larger distribution shift, ImageNet pre-trained models was
 147 found to offer limited benefits for medical imaging tasks due to large distribution shifts induced by
 148 fundamental differences in data sizes, features, and task specifications; that is, lightweight models
 149 perform comparably to standard architectures [44]. To make matters worse, ImageNet performance
 150 is less correlated with and less predictive of downstream performance on diverse tasks beyond
 151 classification such as object detection, few-shot classification, and segmentation [13]. On top of it all,
 152 when ImageNet is extended with a perturbed temporal dimension, models performance significantly
 153 worsen [52].

154 **On the Usability of Benchmarks:** Beyond ImageNet, VTAB introduced a benchmark with a wider
 155 diversity of tasks and domains [70]. Nevertheless, it does not offer task and domain shifts offered
 156 in GATE, such as medical segmentation and video classification and regression that are known to
 157 be ill-measured and gauged by ImageNet alone [44, 52]. That said, VTAB offers satellite imaging
 158 and 3D tasks which GATE does not. Nevertheless, GATE as a software framework was optimized to
 159 minimise usage friction, to take no more than 12 GPU hours on our smallest tier, and, to only require
 160 approximately 1 hour of adding the new encoder and wrapping it into GATE wrappers for GATE to be
 161 able to go away and take care of everything, including dataset downloading, task adapter integration
 162 and full train/val and test cycles with logging of various key metrics. VTAB, in our experience,
 163 requires a lot more manual work in getting the datasets, and integrating new models to be adapted.
 164 Similarly, VLMBench [73] and WILDS [26] offer more diverse datasets beyond previous work but
 165 neither offer a tiered approach that enables iterative development of models during pre-training, nor
 166 produce extensible and flexible benchmarks that can be easily glued into researchers experimentation
 167 code without friction.

168 **On the Systematic Selection of Benchmarks:** Previous work investigated the properties inherit
 169 in multi-task benchmarks that trade-off diversity and sensitivity where the latter is how robust a
 170 benchmark ranking is to the inclusion of irrelevant models or minute changes in the tasks themselves
 171 [72]. It was found that multi-task benchmark are unstable to irrelevant changes in tasks design.
 172 Nevertheless, this is related to how the benchmark ranks models; whether it compares how model often
 173 ranks higher than another in cardinal benchmarks or if the performance across tasks is averaged to
 174 produce a single rank in cardinal ones. Meanwhile, our benchmark produces fine-grained information
 175 to model performances across diverse tasks rather than producing specific ranking which is delegated
 176 to the user analysis. Another complementary thread of work investigates dynamic benchmarks where
 177 model training and data collection is interleaved to continually challenge model knowledge [53]. To
 178 the best of our knowledge, this is the first work that studies the selection of multi-task, multi-domain
 179 benchmarks that satisfy limited compute budgets while maximizing research signal.

180 In summary, Table 1 shows the desiderata that we believe a good evaluation suite and framework
 181 should have such that they can both offer the community useful signal, and also balance that with
 182 being practical so that people can adopt it.

183 3 EEVEE Methodology

184 EEVEE is our proposed method for automating the selection of Pareto-optimal benchmark subsets.
 185 By analyzing benchmark performance metrics, EEVEE identifies a small, highly informative subset
 186 that maximizes information relative to the entire benchmark pool. This ensures that, as machine
 187 learning benchmark breadth and depth increases, we will always be able to identify and select few that
 188 offer high information about the whole. We strike a balance between providing rich evaluation signals
 189 and maintaining simplicity, reducing computational costs and human efforts required for adopting
 190 new benchmarks. EEVEE enables the production of a tiered evaluation engine accommodating
 191 various computational budgets, fostering an inclusive and accessible research environment, and
 192 improving the quality of insights derived from machine learning research while addressing reluctance

193 towards resource-intensive evaluation processes. This balance between efficiency, simplicity, and
 194 signal richness presents EEVEE’s value proposition for advancing machine learning research.

195 **Working Principle of EEVEE:** EEVEE works by building a *meta-model* over the performance
 196 metrics of models sufficient both in number and diversity, on the full benchmark pool from which we
 197 want to choose our subset. With the term *benchmark* in this paper we refer to a dataset + task
 198 pairs.

199 Formally, given a large benchmark pool $B = \{b_0, b_1, \dots, b_K\}$, where B is the full set of benchmarks,
 200 and b_i are individual benchmarks therein, we have a sufficiently large and diverse pool of model
 201 performance metrics $M = \{m_0^0, m_1^0, \dots, m_K^N\}$. Here, m_i^j is the performance metric of model j on
 202 benchmark b_i . We aim to discover a subset of B of size k . This means k total benchmarks make
 203 up the subset. If we build a meta-model $g(M_{selected}, \theta)$ to predict all of M given only the selected
 204 subset $M_{selected}$, it should minimize the following loss:

$$L_{EEVEE} = MSE(M, g(M_{selected}, \theta)) \quad (1)$$

205 In this equation, MSE is the mean squared error. M represents the full set of performance metrics of
 206 all our models on the full benchmark pool B . The term $g(M_{selected}, \theta)$ represents the predictions of
 207 the meta-model g with parameters θ when it is given the performance metrics of all models from the
 208 selected subset of benchmarks $B_{selected}$, referred to as $M_{selected}$.

209 However, our main focus lies in the selected combination of performance metrics $M_{selected}$ that can
 210 generalize well on previously unseen models. To that end, we must split M into train, validation
 211 and test sets, each consisting of performance metrics acquired from different models (e.g. train
 212 \rightarrow ResNet50, ViT-Base, CLIP, and val \rightarrow ResNext50, DINO, DeiT), and explicitly optimize the
 213 inner loop test loss rather than the training loss, while we use the validation loss to select the best
 214 meta-model for test. Hence the loss we wish to minimize is:

$$L_{EEVEE}^{test} = MSE(M^{test}, g(M_{selected}^{test}, \theta)) \quad (2)$$

215 We need a non-differentiable method for choosing the k benchmarks in $M_{selected}$, since brute
 216 force becomes intractable very quickly, so we employ evolutionary methods to learn the k selected
 217 benchmarks.

218 This results in a bi-level optimization, with an evolutionary method on the outer loop $e(B_{selected})$,
 219 where e is the evolutionary method, and $B_{selected}$ are the benchmarks being selected – or indeed, the
 220 genes being optimized, and a small meta-model parameterized as a neural network $g(\theta)$ that receives
 221 a train/val split from $B_{selected}$ and trains itself to do the task described in Equation 1, after which
 222 process it is scored using the val set using the loss in Equation 2. Then, once a given candidate of
 223 benchmarks $B_{selected}$ is scored, in this way, the outer loop performs a tournament selection where
 224 only the top 50 candidates are preserved and mutated by removing one benchmark at random, and
 225 adding another at random. Each winning candidate mutates into 10 children, and the parent is
 226 also preserved in the gene pool, producing a gene pool with 550 candidates for every cycle. At
 227 initialization, we sample 1000 random combinations. We have found that 1000 is a good starting
 228 population that is both tractable to score and facilitates the necessary diversity that enables limited
 229 variation in results across several runs, showcasing convergent behaviour. diversity that our results
 230 across runs have little variation from one another, pointing to a convergent behaviour. We include full
 231 pseudocode showcasing all the details related to how we performed EEVEE for our experiments in
 232 Algorithm 1, 2 and 3 in Figure 1

233 Applying EEVEE on Model Encoder Generalization

234 **Why Model Encoder Evaluation?** A common practice across machine learning applications involves
 235 augmenting general model encoders with task-oriented heads. The adaption of this paradigm can
 236 be attributed to the computational efficiency associated with training model encoders, over more
 237 expensive setups. Much of computer vision, as well as vision to text search and retrieval happen using
 238 model encoders. Similarly, various applications requiring translation from one domain/modality/task
 239 to another require an encoder of some sort. Even the “decoder-only” LLM models that have
 240 demonstrated incredible capabilities in the last few years, internally can be seen as a series of
 241 representation encoders, a series of refinement before they reach the decoding stage.

242 Multi-modal early fusion is another
 243 topic closely related with model encod-
 244 ers – as research in early fusion
 245 can be done most efficiently when try-
 246 ing to learn data encoders rather than
 247 a full encoder-decoder, or decoder-
 248 only models. World model research
 249 in multi-modal dimensions can also
 250 take place most efficiently within a
 251 model-encoder context. Recent works
 252 like I/VJEPa [2] for example have
 253 paved the way for self-supervised
 254 learning which functions using model
 255 encoders, and has been demonstrated
 256 to be more efficient and more gener-
 257 alizable than full pixel decoding vari-
 258 ants.

259 Furthermore, model encoder evalua-
 260 tion has been quite diffused in the past
 261 few years, with new benchmarks be-
 262 ing produced in every facet of the machine learning field. Nonetheless, most of those lacked in some
 263 key quality: they were either simply too complex to use efficiently, requiring too much compute, or,
 264 more often than the others, missing a unifying software framework that can easily, in a user-conscious
 265 way, and a principled stance towards high readability, maintainability and hackability.

266 **The goal of focusing on Model Encoder Eval-**
 267 **uation:** By applying EEVEE to search for a
 268 pareto-optimal set of benchmarks, *and* packag-
 269 ing it up in a unified framework that is built for
 270 the researcher in mind from the ground up, one
 271 which offers out of the box automated down-
 272 loading, pipeline building, task adapters, and a
 273 very mature training and eval loop. Within this
 274 framework, we facilitate, all relevant logging in-
 275 formation, including key training and eval met-
 276 rics, rich gradient information, power and com-
 277 putational information, as well as visualizations
 278 where relevant. Finally, we support easy switching of model encoders, no matter what source modality
 279 they come from – our framework dubbed *GATE* is a one stop shop for ones model representation
 280 research needs, both during research, debugging, as well as at the evaluation phase.

281 *GATE* comes in three tiers *small*, *base* and *big-GATE*. Each having 8, 15 and 21 benchmarks within it,
 282 and targetted towards 12/24 and 36 GPU hours on a A100 40GB. We hope that by making it very easy
 283 for the end user and offering such rich signal for machine learning research, many researchers will
 284 choose to use *GATE*, to enhance their research signal, whilst keeping the compute budgets relatively
 285 feasible.

286 **Preparations: Choosing Models, Benchmarks and Adaptation Processes:** EEVEE will yield
 287 better results if the space of models, benchmarks and adaptation processes we use is diverse, but also
 288 thorough in numbers. **A. Adaptation Process** We wanted *GATE* to cover multiple domains, tasks
 289 and modalities when shifting from the source to the target setting. For that reason we decided that if
 290 a model encoder has an input layer that does not fit the target modality, we simply remove that input
 291 layer and replace it with a relevant ViT-like patchification [12] followed by a linear combination for
 292 each patch. For tasks where we have text, we would tokenize the text using BPE [51], and for tasks
 293 where we have video we would use the model encoder on each image, to acquire an image-level
 294 vector representation, and then follow that up with a simple 4 layer transformer that receives a
 295 sequence of image-vector tokens, to produce a video-level embedding, on top of which we apply the
 296 task-specific head at hand. The task-adapters we used leaned on established methods, and where
 297 possible we just used a transformer head, which includes segmentation, relational reasoning and
 298 video classification, with everything just using a linear head, full details available at 14. After these

Algorithm 1 Scoring

Require: Performance metrics M , Input metrics M_{selected} ,
 Epochs $E = 20$, Hidden dimension $d_{\text{hidden}} = 100$,
 Learning rate $\alpha = 0.01$, Weight decay $\lambda = 0.01$, Opti-
 mizer type $\omega = \text{"AdamW"}$
Ensure: Evaluation score $\text{mean}(\text{scores})$
 1: Convert data to tensors $x = M_{\text{selected}}$ and $y = M$
 2: Normalize x and y
 3: Initialize ShuffleSplit cross-validation kf
 4: Initialize empty list scores
 5: for each train, val split in kf do
 6: Divide x into x_{train} and x_{val} ; y into y_{train} and y_{val}
 7: Build meta-model $g(\theta)$ with hidden dimension d_{hidden}
 8: Train $g(\theta)$ on x_{train} and y_{train} for E epochs with learn-
 ing rate α , weight decay λ , and optimizer ω
 9: Predict $y_{\text{pred}} = g(x_{\text{val}}, \theta)$
 10: Compute mean squared error score =
 MSE($y_{\text{pred}}, y_{\text{val}}$)
 11: Append score to scores
 12: end for
 13: return $\text{mean}(\text{scores})$

Algorithm 2 Mutation

Require: $B_{\text{selected}} \subset B$, $B = \{b_1, b_2, \dots, b_k\}$
Ensure: New B_{selected}
 1: Select $b_{\text{remove}} \in B_{\text{selected}}$
 2: Select $b_{\text{add}} \in B$
 3: while $b_{\text{add}} \in B_{\text{selected}}$ do
 4: Select another $b_{\text{add}} \in B$
 5: end while
 6: Create B_{selected} by replacing b_{remove} with b_{add}
 7: return B_{selected}

Algorithm 3 Evolution

Require: Performance metrics $M = \{m_1^1, m_1^2, \dots, m_k^N\}$,
 Benchmark set B , Combination size k , Number of win-
 ners W , Number of children per winner C , Number
 of generations G , Initial combinations size I , Training
 epochs E , Hidden dimension $d_{\text{hidden}} = 100$, Learning
 rate $\alpha = 0.01$, Weight decay $\lambda = 0.01$, Optimizer type
 $\omega = \text{"AdamW"}$
Ensure: Evolved benchmark combinations B_{winners}
 1: Initialize initial combinations B_{initial} with I random sam-
 ples from B of size k
 2: Evaluate performance of B_{initial} using SCOR-
 ING($M, B_{\text{initial}}, E, d_{\text{hidden}}, \alpha, \lambda, \omega$) and store scores in
 S
 3: Select top W combinations from S as B_{winners}
 4: for generation $g = 1$ to G do
 5: Initialize a new set of combinations B_{new}
 6: for each combination $B_{\text{selected}} \in B_{\text{winners}}$ do
 7: Add B_{selected} to B_{new}
 8: for each child $c = 1$ to C do
 9: Mutate B_{selected} using MUTATION(B_{selected}, B)
 to create a new combination B_{selected}
 10: Add B_{selected} to B_{new}
 11: end for
 12: end for
 13: Evaluate performance of B_{new} using SCOR-
 ING($M, B_{\text{new}}, E, d_{\text{hidden}}, \alpha, \lambda, \omega$) and store scores
 in S
 14: Select top W combinations from S as B_{winners}
 15: end for
 16: return B_{winners}

Figure 1: (a) EEVEE Scoring algorithm, Mutation algorithm, and (b) Evolution algorithm.

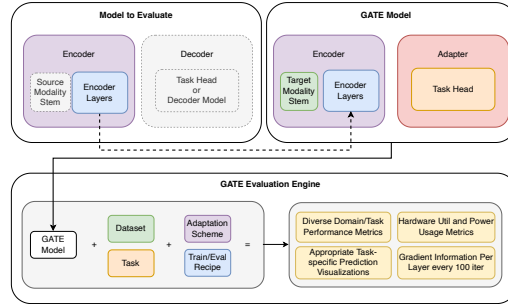


Figure 2: GATE Framework Pipeline

299 modifications, described in Figure 2, we use a fine tuning scheme – this decision was informed by
300 preliminary experiments on both full fine tuning and linear probe with a frozen backbone, in which
301 we found that there was a clear superiority of fine tuning over linear probing for the benchmarks we
302 chose in our pool. Full details of these preliminary experiments can be found in Appendix 8.1. In our
303 preliminary experiments we were able to identify three recipes, one for ConvNet-style architectures,
304 one for ViT-style architectures and one for Hybrid architectures such as ConvNext and ResNext that
305 worked well for all tasks, details in 8.1.

306 **B. Model Pool** We wanted the space of models used to cover many important pretraining schemes,
307 architectures, and source modalities. The details of these choices are provided next: **1. Pretraining**
308 **Task and Dataset Variation:** With a consistent architecture, models were subjected to various
309 pretraining tasks and datasets. Model instances representing this category include CLIPViT [43],
310 ConvNextV2 [35], Siglip, FlexViT [7], LaionViT, ImageNet1K ViT [11] with Random Aug-
311 ment, SAM-ViT, DiNoViT, EfficientFormerV2 [32] and DeiT3 [59]. Further to these, we include
312 models initialized from scratch, specifically, ViT, ResNet50 [18], FlexViT, EfficientNetV2 [57],
313 and then fine-tuned on the GATE tasks. **2. Architectural Variation:** We explored models having the
314 same pretraining dataset (ImageNet), but differing in their architecture. This group encompassed a
315 mix of standard CNN models such as EffNetV2, ResNet50, ResNext50 [67], ConvNextV2_Base
316 [35] and transformer-based models like EfficientFormer [32] and FlexViT [7]. **3. Modality**
317 **and Dataset Variation:** This axis comprised models trained on modalities other than vision such
318 as Whisper, coming from an audio to text task and Bert [10], Bart [31] and MpNet [55] coming
319 from various text-based tasks. These models had their original input processing systems replaced by
320 a Vision Transformer style embedding and were subsequently fine-tuned on the GATE tasks. A more
321 comprehensive account of these models, including their selection rationale and unique characteristics,
322 is provided in the Appendix Section 13.

323 **C. Benchmark Pool** The benchmark pool, detailed in the Appendix, includes Image Classification
324 (ImageNet1k [9], CIFAR100 [28], Places365 [74], Food101 [36], HappyWhale [17]), Few Shot
325 Image Classification (Aircraft [37], Fungi [50], MiniImageNet [62], CUB200 [63], Describable
326 Features [69]), Zero Shot Text-Image Classification (Flickr30K [41], New Yorker Caption Context
327 [20], Winoground [58]), Visual Relational Reasoning (CLEVR [23], CLEVRMath [34]), Image
328 Semantic Segmentation (ADE20K [75], COCO10K [33], COCO164K [33], NYU-Depth-v2 [54],
329 PascalContext [38], Cityscapes [8]), Medical Image Classification (Chexpert [21], Diabetic Retinopa-
330 thy [16], HAM10000 [60]), Medical Segmentation (ACDC [5]), Video Classification (HMDB51 [30],
331 UCF-101 [56], Kinetics400 [24]) and Video Regression (iWildcam [4]).

332 **Producing Diverse Model Performance Metrics:** We apply our adaptation process on each and
333 every model chosen, on every benchmark in the benchmark pool. To acquire test results we ensemble
334 by averaging logits of the top 1, 3 and 5 validation models to produce three separate ensemble results.

335 **D. Experimental Approach** We wanted our research environment to reflect the end user, so we
336 can properly understand their needs, and to offer a *pragmatic* experimental setup of in-the-wild
337 researchers with little time to hyperparameter optimize, and which have to make decisions on small
338 amounts of preliminary experiments – someone choosing a model encoder off the shelf and adapting it
339 to downstream setting. For that reason, we kept any hyperparameter tuning, or human attention when
340 it came to specific models to a minimum. Instead, we relied on existing good recipes, and did some
341 preliminary experiments as explained in detail in 8.1. Briefly, we discovered specific adjustments
342 for each architecture type: for Convolutional Architectures, we used AdamW with a learning rate of
343 $1e-3$, and $6e-4$ for segmentation tasks; for Vision Transformer Architectures, AdamW with a learning
344 rate of $1e-5$; and for Convolutional + Transformer Hybrid Architectures, AdamW with a learning rate
345 of $2e-5$. A plateau learning rate scheduler was configured with parameters like mode "min", factor
346 0.5, patience 1000, and threshold $1e-4$, allowing models to effectively choose their own schedules
347 based on their learning progress. This adaptive scheduling facilitated “good enough” learning rates
348 and enhanced performance across different architectures.

349 4 Results

350 **Single Benchmark Predictiveness:** As demonstrated in Figure 3, using EEVEE we quantified the
351 predictive power of each benchmark **on its own**, when not in a combination with others. We have
352 found that ADE20K, Flickr30K, and the New York Caption Competition lead in their predictive
353 power, with few-shot tasks, and relational reasoning, being very close to the best in predictive power.
354 ImageNet1K sits squarely in the middle of the competition. Furthermore, some of the most “novel”

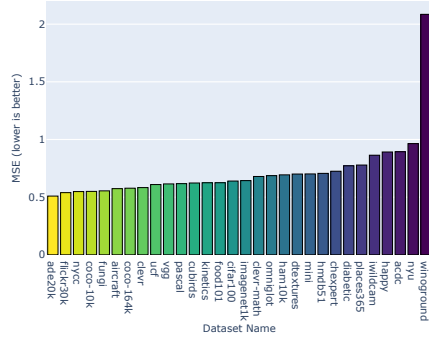
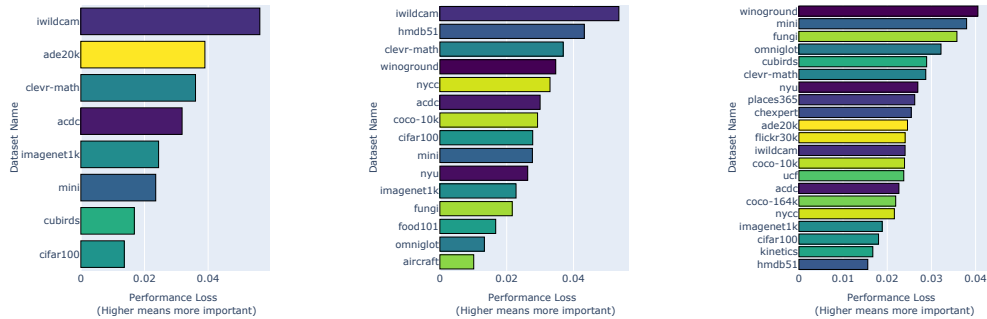


Figure 3: The EEVEE MSE Loss ($k=1$) shows "predictiveness over the whole," with lower values being better. Benchmarks like iWildcam, HappyWhale, and WinoGround test unique capabilities and may not predict all tasks, yet EEVEE often includes at least two of these in its top combinations along with a "natural-image representative" such as CIFAR100, ADE20K or Flickr30K.



(a) Small-GATE ($k=8$, 12 GPU hour) tier (b) Small-GATE ($k=15$, 24 GPU hour) tier (c) Small-GATE ($k=21$, 24 GPU hour) tier

Figure 4: Degradation of predictive power when a given benchmark is removed and the meta-model trained from scratch, for different GATE tiers.

355 benchmarks like iwildcam, happy whale, ACDC, NYU and Winoground are the least predictive tasks,
 356 Winoground being magnitudes less predictive. We argue that this is mainly due to the tasks being
 357 "harder", and our models being less designed for those. The results in WinoGround were barely better
 358 than chance for example. However, when once we move to combinations of benchmarks, these 'less'
 359 predictive benchmarks become key contributors to better predictive power, as they represent edge
 360 cases, as can be seen in Figures 6g 7c, 7i, where these have the highest importance when removed
 361 from a given set.

362 Predictiveness of Discovered Combinations In

363 Figure 5, we can see how the top-50 performing
 364 candidate combinations perform as we vary the
 365 number of benchmarks per combination from
 366 1 to 26. We can see that there is a point of
 367 diminishing returns around the $k = 8$ point, after
 368 which there appears to be some "overfitting"
 369 occurring. We verified that the overfitting was a
 370 result of having a small sample number of 20
 371 models, to train, val and test our meta-models
 372 with, as well as the 2-layer MLP we used to
 373 model Few-to-All metric predictions. We tried
 374 our level best to find the best architecture and
 375 regularization schemes for our meta-model, and
 376 this was the best we could do given available
 377 compute and (human) time. We chose 8, 15,
 378 and 21 as the combination-threshold to make
 379 our packs out of as they satisfied the computa-
 380 tional budgets we set for ourselves, and they
 381 have very diverse and predictive tasks, as can
 382 be seen in Figures 6g 7c, 7i. For full details
 383 on all the discovered top-k combinations please
 384 look at Appendix Section 16.1. **Best Models based on GATE:** As can be seen in Table 2, or the Appendix extended Table 3, the best overall models are ConvNextV2, SigLIP and CLIP in that order, with SigLIP and CLIP often exchanging ranks between themselves. However, it is worth noting that EfficientNetV2 demonstrated exceptional performance/compute across all tasks, and even outperformed all models in many medical tasks. Finally, ConvNet based models, and particularly ResNext50 seem to have done exceptionally well in the edge-case scenarios of ACDC, Happy Whale Individual identification, and

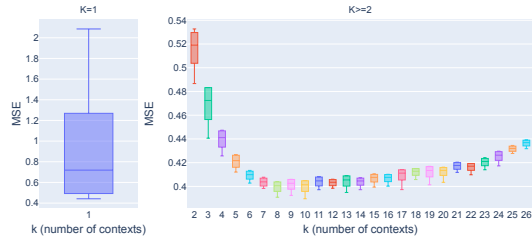


Figure 5: Performance of Models build with K -best datasets: We do a search over the space of all k for EEVEE and box plot the population summary statistics of the top 50 combination candidates.

References

- 406
- 407 [1] Stanislaw Antol, Aishwarya Agrawal, Jiasen Lu, Margaret Mitchell, Dhruv Batra,
408 C Lawrence Zitnick, and Devi Parikh. Vqa: Visual question answering. In *Proceedings*
409 *of the IEEE International Conference on Computer Vision (ICCV)*, pages 2425–2433, 2015.
- 410 [2] Adrien Bardes, Quentin Garrido, Jean Ponce, Xinlei Chen, Michael Rabbat, Yann LeCun,
411 Mahmoud Assran, and Nicolas Ballas. Revisiting feature prediction for learning visual repre-
412 sentations from video, 2024.
- 413 [3] Loic Barrault, Ondrej Bojar, Marta R Costa-jussa, Christian Federmann, Mark Fishel, Yvette
414 Graham, Barry Haddow, Matthias Huck, Philipp Koehn, Shervin Malmasi, Christof Monz, et al.
415 Findings of the 2019 conference on machine translation (wmt19). In *Proceedings of the Fourth*
416 *Conference on Machine Translation (Volume 2: Shared Task Papers, Day 1)*, pages 1–61, 2019.
- 417 [4] Sara Beery, Grant Van Horn, and Pietro Perona. The iwildcam 2018 challenge dataset. In
418 *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*
419 *Workshops*, pages 54–60, 2018.
- 420 [5] Olivier Bernard, Alain Lalonde, Caio Zotti, Florence Cervenansky, Xin Yang, Pheng-Ann
421 Heng, Ismail Cetin, Karim Lekadir, Oscar Camara, Miguel A Gonzalez Ballester, et al. Deep
422 learning techniques for automatic mri cardiac multi-structures segmentation and diagnosis: Is
423 the problem solved? *IEEE Transactions on Medical Imaging*, 37(11):2514–2525, 2018.
- 424 [6] Lucas Beyer, Olivier J. Hénaff, Alexander Kolesnikov, Xiaohua Zhai, and Aäron van den Oord.
425 Are we done with imagenet?, 2020.
- 426 [7] Lucas Beyer, Pavel Izmailov, Alexander Kolesnikov, Mathilde Caron, Simon Kornblith, Xiaohua
427 Zhai, Matthias Minderer, Michael Tschannen, Ibrahim M. Alabdulmohsin, and Filip Pavetic.
428 Flexivit: One model for all patch sizes. *2023 IEEE/CVF Conference on Computer Vision and*
429 *Pattern Recognition (CVPR)*, pages 14496–14506, 2022.
- 430 [8] Marius Cordts, Mohamed Omran, Sebastian Ramos, Tobias Rehfeld, Markus Enzweiler, Rodrigo
431 Benenson, Uwe Franke, Stefan Roth, and Bernt Schiele. The cityscapes dataset for semantic
432 urban scene understanding. In *Proceedings of the IEEE Conference on Computer Vision and*
433 *Pattern Recognition*, pages 3213–3223, 2016.
- 434 [9] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Li Kai, and Li Fei-Fei. Imagenet: A large-
435 scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern*
436 *recognition*, pages 248–255. Ieee, 2009.
- 437 [10] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of
438 deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*,
439 2018.
- 440 [11] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai,
441 Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly,
442 Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image
443 recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020.
- 444 [12] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai,
445 Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly,
446 Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image
447 recognition at scale. In *International Conference on Learning Representations (ICLR)*, 2021.
- 448 [13] Linus Ericsson, Henry Gouk, and Timothy M. Hospedales. How well do self-supervised
449 models transfer? In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern*
450 *Recognition (CVPR)*, pages 5414–5423, June 2021.
- 451 [14] Mark Everingham, Luc Van Gool, Christopher KI Williams, John Winn, and Andrew Zisserman.
452 The pascal visual object classes (voc) challenge. *International Journal of Computer Vision*,
453 88(2):303–338, 2010.

- 454 [15] John S Garofolo, Lori F Lamel, William M Fisher, Jonathan G Fiscus, David S Pallett, and
455 Nancy L Dahlgren. Timit acoustic-phonetic continuous speech corpus ldc93s1, 1993.
- 456 [16] Varun Gulshan, Lily Peng, Marc Coram, Michael C Stumpe, Derek Wu, Arunachalam
457 Narayanaswamy, Subhashini Venugopalan, Kasumi Widner, Travis Madams, Jorge Cuadros,
458 et al. Development and validation of a deep learning algorithm for detection of diabetic
459 retinopathy in retinal fundus photographs. *JAMA*, 316(22):2402–2410, 2016.
- 460 [17] Happywhale. Happywhale - whale and dolphin identification challenge, 2022.
- 461 [18] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image
462 recognition. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*,
463 pages 770–778. IEEE, 2016.
- 464 [19] Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Mantas He, Dawn
465 Song, and Jacob Steinhardt. Measuring massive multitask language understanding. *arXiv*
466 *preprint arXiv:2009.03300*, 2020.
- 467 [20] Jack Hessel. New yorker caption contest corpus, 2023.
- 468 [21] Jeremy Irvin, Pranav Rajpurkar, Michael Ko, Yifan Yu, Silvana Ciurea-Ilcus, Chris Chute,
469 Henrik Marklund, Behzad Haghgoo, Robyn Ball, Katie Shpanskaya, et al. Chexpert: A large
470 chest radiograph dataset with uncertainty labels and expert comparison. *Proceedings of the*
471 *AAAI Conference on Artificial Intelligence*, 33:590–597, 2019.
- 472 [22] Nihues Jan et al. Iwslt 2017: Proceedings of the 14th international workshop on spoken
473 language translation. 2017.
- 474 [23] Justin Johnson, Bharath Hariharan, Laurens van der Maaten, Li Fei-Fei, C Lawrence Zitnick,
475 and Ross Girshick. Clevr: A diagnostic dataset for compositional language and elementary
476 visual reasoning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern*
477 *Recognition*, pages 2901–2910, 2017.
- 478 [24] Will Kay, Joao Carreira, Karen Simonyan, Brian Zhang, Chloe Hillier, Sudheendra Vijaya-
479 narasimhan, Fabio Viola, Tim Green, Trevor Back, Paul Natsev, et al. The kinetics human
480 action video dataset. In *arXiv preprint arXiv:1705.06950*, 2017.
- 481 [25] Philipp Koehn. Europarl: A parallel corpus for statistical machine translation. *MT summit*,
482 5:79–86, 2005.
- 483 [26] Pang Wei Koh, Shiori Sagawa, Henrik Marklund, Sang Michael Xie, Marvin Zhang, Akshay
484 Balsubramani, Weihua Hu, Michihiro Yasunaga, Richard Phillips, Irena Gao, et al. Wilds: A
485 benchmark of in-the-wild distribution shifts. In *International Conference on Machine Learning*,
486 2021.
- 487 [27] Simon Kornblith, Jonathon Shlens, and Quoc V. Le. Do better imagenet models transfer better?
488 In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*
489 *(CVPR)*, June 2019.
- 490 [28] Alex Krizhevsky and Geoffrey Hinton. Learning multiple layers of features from tiny images.
491 Technical Report UTML TR 2009, University of Toronto, Toronto, Ontario, Canada, 2009.
- 492 [29] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep
493 convolutional neural networks. *Advances in neural information processing systems*, 25:1097–
494 1105, 2012.
- 495 [30] Hildegard Kuehne, Hueihan Jhuang, Estibaliz Garrote, Tomaso Poggio, and Thomas Serre.
496 Hmdb: A large video database for human motion recognition. In *2011 International Conference*
497 *on Computer Vision (ICCV)*, pages 2556–2563. IEEE, 2011.
- 498 [31] Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed,
499 Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. Bart: Denoising sequence-to-sequence
500 pre-training for natural language generation, translation, and comprehension. *arXiv preprint*
501 *arXiv:1910.13461*, 2020.

- 502 [32] Xiuyu Li, Yutong Yuan, Shu Chen, Martin Danelljan, Radu Timofte, and Luc Van Gool.
503 Efficientformer: Vision transformers at mobilenet speed. *arXiv preprint arXiv:2206.01191*,
504 2022.
- 505 [33] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr
506 Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *European*
507 *Conference on Computer Vision*, pages 740–755, 2014.
- 508 [34] Adam Dahlgren Lindström and Savitha Sam Abraham. Clevr-math: A dataset for compositional
509 language, visual and mathematical reasoning. *ArXiv*, abs/2208.05358, 2022.
- 510 [35] Zhuang Liu, Hanzi Mao, Chao-Yuan Wu, Christoph Feichtenhofer, Trevor Darrell, and Saining
511 Xie. A convnet for the 2020s. *arXiv preprint arXiv:2201.03545*, 2022.
- 512 [36] Matthieu Guillaumin Lukas Bossard and Luc Van Gool. Food-101 – mining discriminative
513 components with random forests. In *European Conference on Computer Vision*, 2014.
- 514 [37] Subhransu Maji, Esa Rahtu, Juho Kannala, Matthew Blaschko, and Andrea Vedaldi. Fine-
515 grained visual classification of aircraft. In *2013 IEEE Conference on Computer Vision and*
516 *Pattern Recognition (CVPR)*, pages 554–561. IEEE, 2013.
- 517 [38] Roozbeh Mottaghi, Xiaobai Chen, Xiaofeng Liu, Nam-Gyu Cho, Seong-Whan Lee, Sanja
518 Fidler, Raquel Urtasun, and Alan Yuille. The role of context for object detection and semantic
519 segmentation in the wild. In *Proceedings of the IEEE Conference on Computer Vision and*
520 *Pattern Recognition (CVPR)*, pages 891–898, 2014.
- 521 [39] Vassil Panayotov, Guoguo Chen, Daniel Povey, and Sanjeev Khudanpur. Librispeech: An
522 asr corpus based on public domain audio books. In *2015 IEEE International Conference on*
523 *Acoustics, Speech and Signal Processing (ICASSP)*, pages 5206–5210, 2015.
- 524 [40] Denis Paperno, Germán Kruszewski, Angeliki Lazaridou, Ngoc Quan Pham, Raffaella Bernardi,
525 Sandro Pezzelle, Marco Baroni, Gemma Boleda, and Raquel Fernández. The lambada dataset:
526 Word prediction requiring a broad discourse context. In *Proceedings of the 54th Annual Meeting*
527 *of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1525–1534,
528 2016.
- 529 [41] Bryan Plummer, Liwei Wang, Chris Cervantes, Juan Caicedo, Julia Hockenmaier, and Svetlana
530 Lazebnik. Flickr30k entities: Collecting region-to-phrase correspondences for grounded image
531 descriptions. *arXiv preprint arXiv:1505.04870*, 2015.
- 532 [42] Bryan A Plummer, Liwei Wang, Christopher M Cervantes, Juan C Caicedo, Julia Hockenmaier,
533 and Svetlana Lazebnik. Flickr30k entities: Collecting region-to-phrase correspondences for
534 richer image-to-sentence models. In *Proceedings of the IEEE International Conference on*
535 *Computer Vision*, pages 2641–2649, 2015.
- 536 [43] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,
537 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual
538 models from natural language supervision. In *International conference on machine learning*,
539 pages 8748–8763. PMLR, 2021.
- 540 [44] Maithra Raghu, Chiyuan Zhang, Jon Kleinberg, and Samy Bengio. Transfusion: Understanding
541 transfer learning for medical imaging. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-
542 Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems*,
543 volume 32. Curran Associates, Inc., 2019.
- 544 [45] Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. Squad: 100,000+ questions
545 for machine comprehension of text. In *Proceedings of the 2016 Conference on Empirical*
546 *Methods in Natural Language Processing*, pages 2383–2392, 2016.
- 547 [46] Benjamin Recht, Rebecca Roelofs, Ludwig Schmidt, and Vaishal Shankar. Do cifar-10
548 classifiers generalize to cifar-10?, 2018.

- 549 [47] Benjamin Recht, Rebecca Roelofs, Ludwig Schmidt, and Vaishaal Shankar. Do ImageNet
550 classifiers generalize to ImageNet? In *Proceedings of the 36th International Conference on*
551 *Machine Learning*, Proceedings of Machine Learning Research. PMLR, 2019.
- 552 [48] Erik F Sang and Fien De Meulder. Introduction to the conll-2003 shared task: Language-
553 independent named entity recognition. In *Proceedings of the seventh conference on Natural*
554 *language learning at HLT-NAACL 2003*, pages 142–147, 2003.
- 555 [49] Adam Santoro, David Raposo, David G.T. Barrett, Mateusz Malinowski, Razvan Pascanu, Peter
556 Battaglia, and Timothy Lillicrap. A simple neural network module for relational reasoning. In
557 *Advances in Neural Information Processing Systems (NeurIPS)*, pages 4967–4976, 2017.
- 558 [50] Dirk Schroeder, Yin Cui, Yang Chai, Daniel Kristensen, Evangelos Kalogerakis, and Serge
559 Belongie. The fgvcx fungi classification challenge. In *CVPR Workshops*, 2018.
- 560 [51] Rico Sennrich, Barry Haddow, and Alexandra Birch. Neural machine translation of rare
561 words with subword units. In *Proceedings of the 54th Annual Meeting of the Association for*
562 *Computational Linguistics (Volume 1: Long Papers)*, pages 1715–1725, 2016.
- 563 [52] Vaishaal Shankar, Achal Dave, Rebecca Roelofs, Deva Ramanan, Benjamin Recht, and Ludwig
564 Schmidt. Do image classifiers generalize across time? In *Proceedings of the IEEE/CVF*
565 *International Conference on Computer Vision (ICCV)*, pages 9661–9669, October 2021.
- 566 [53] Ali Shirali, Rediet Abebe, and Moritz Hardt. A theory of dynamic benchmarks, 2023.
- 567 [54] Nathan Silberman, Derek Hoiem, Pushmeet Kohli, and Rob Fergus. Indoor segmentation and
568 support inference from rgb-d images. In *Proceedings of the European Conference on Computer*
569 *Vision (ECCV)*, pages 746–760, 2012.
- 570 [55] Kaitao Song, Xu Tan, Tao Qin, Jianfeng Lu, and Tie-Yan Liu. MpNet: Masked and permuted
571 pre-training for language understanding. *arXiv preprint arXiv:2004.09297*, 2020.
- 572 [56] Khurram Soomro, Amir Roshan Zamir, and Mubarak Shah. Ucf101: A dataset of 101 human
573 actions classes from videos in the wild. In *arXiv preprint arXiv:1212.0402*, 2012.
- 574 [57] Mingxing Tan and Quoc V Le. Efficientnetv2: Smaller models and faster training. *arXiv*
575 *preprint arXiv:2104.00298*, 2021.
- 576 [58] Tristan Thrush, Hongyu Jiang, Goutham Prasad, and Jacob Andreas. Winoground: Prob-
577 ing vision and language models for visio-linguistic compositionality. *arXiv preprint*
578 *arXiv:2204.03162*, 2022.
- 579 [59] Hugo Touvron, Piotr Bojanowski, Mathilde Caron, Matthieu Cord, Alaaeldin El-Nouby,
580 Edouard Grave, Armand Joulin, Gabriel Synnaeve, and Jakob Verbeek. Deit iii: Revenge
581 of the vit. *arXiv preprint arXiv:2204.07118*, 2022.
- 582 [60] Philipp Tschandl, Cliff Rosendahl, and Harald Kittler. The ham10000 dataset, a large collection
583 of multi-source dermatoscopic images of common pigmented skin lesions. *Scientific Data*,
584 5:180161, 2018.
- 585 [61] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez,
586 Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in Neural Informa-*
587 *tion Processing Systems (NeurIPS)*, pages 5998–6008, 2017.
- 588 [62] Oriol Vinyals, Charles Blundell, Timothy Lillicrap, Koray Kavukcuoglu, and Daan Wierstra.
589 Matching networks for one shot learning. In *Advances in Neural Information Processing*
590 *Systems (NeurIPS)*, pages 3630–3638, 2016.
- 591 [63] Catherine Wah, Steve Branson, Peter Welinder, Pietro Perona, and Serge Belongie. The caltech-
592 ucsd birds-200-2011 dataset. Technical Report CNS-TR-2011-001, California Institute of
593 Technology, 2011.

- 594 [64] Alex Wang, Yada Pruksachatkun, Nikita Nangia, Amanpreet Singh, Julian Michael, Felix Hill,
595 Omer Levy, and Samuel R Bowman. Superglue: A stickier benchmark for general-purpose
596 language understanding systems. In *Proceedings of the 33rd International Conference on*
597 *Neural Information Processing Systems*, pages 3266–3280, 2019.
- 598 [65] Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R Bowman.
599 Glue: A multi-task benchmark and analysis platform for natural language understanding. In
600 *Proceedings of the 2018 EMNLP Workshop BlackboxNLP: Analyzing and Interpreting Neural*
601 *Networks for NLP*, pages 353–355, 2018.
- 602 [66] Jason Weston, Antoine Bordes, Sumit Chopra, Alexander M Rush, Bart van Merriënboer,
603 Armand Joulin, and Tomas Mikolov. Towards ai-complete question answering: A set of
604 prerequisite toy tasks. *arXiv preprint arXiv:1502.05698*, 2015.
- 605 [67] Saining Xie, Ross Girshick, Piotr Dollár, Zhuowen Tu, and Kaiming He. Aggregated residual
606 transformations for deep neural networks. In *Proceedings of the IEEE Conference on Computer*
607 *Vision and Pattern Recognition (CVPR)*, pages 1492–1500. IEEE, 2017.
- 608 [68] Junichi Yamagishi, Christophe Veaux, and Kirsten MacDonald. Cstr vctk corpus: English
609 multi-speaker corpus for cstr voice cloning toolkit (version 0.92). *University of Edinburgh. The*
610 *Centre for Speech Technology Research (CSTR)*, 2019.
- 611 [69] Fisher Yu, Yinda Zhang, Shuran Song, Ari Seff, and Jianxiong Xiao. Fine-grained visual
612 comparisons with local learning. In *2014 IEEE Conference on Computer Vision and Pattern*
613 *Recognition (CVPR)*, pages 192–199. IEEE, 2014.
- 614 [70] Xiaohua Zhai, Joan Puigcerver, Alexander Kolesnikov, Pierre Ruysen, Carlos Riquelme, Mario
615 Lucic, Josip Djolonga, Andre Susano Pinto, Maxim Neumann, Alexey Dosovitskiy, et al. A
616 large-scale study of representation learning with the visual task adaptation benchmark. In
617 *International Conference on Learning Representations*, 2020.
- 618 [71] Chi Zhang, Feng Gao, Baoxiong Jia, Song-Chun Zhu, and Yixin Zhu. Raven: A dataset for
619 relational and analogical visual reasoning. In *Proceedings of the IEEE Conference on Computer*
620 *Vision and Pattern Recognition*, pages 5317–5327, 2019.
- 621 [72] Guanhua Zhang and Moritz Hardt. Inherent trade-offs between diversity and stability in
622 multi-task benchmarks, 2024.
- 623 [73] Kaizhi Zheng, Xiaotong Chen, Odest Chadwicke Jenkins, and Xin Eric Wang. VLMbench: A
624 compositional benchmark for vision-and-language manipulation. In *Thirty-sixth Conference on*
625 *Neural Information Processing Systems (NeurIPS) Datasets and Benchmarks*, 2022.
- 626 [74] Bolei Zhou, Agata Lapedriza, Aditya Khosla, Aude Oliva, and Antonio Torralba. Places: A 10
627 million image database for scene recognition. In *IEEE Transactions on Pattern Analysis and*
628 *Machine Intelligence*, volume 40, pages 1452–1464. IEEE, 2017.
- 629 [75] Bolei Zhou, Hang Zhao, Xavier Puig, Sanja Fidler, Adela Barriuso, and Antonio Torralba.
630 Scene parsing through ade20k dataset. In *Proceedings of the IEEE Conference on Computer*
631 *Vision and Pattern Recognition (CVPR)*, pages 5122–5130, 2017.
- 632 [76] Yuke Zhu, Olaf Groth, Michael S Bernstein, and Li Fei-Fei. Visual7w: Grounded question
633 answering in images. In *Proceedings of the IEEE Conference on Computer Vision and Pattern*
634 *Recognition*, pages 4995–5004, 2016.

635 6 End-user Guidelines

636 For an end-user to use GATE, they need to:

- 637 1. Install the GATE framework python package, as described in the Github repo’s readme page.
- 638 2. Choose a path for implementing the new foundation model encoder they wish to evaluate.
639 This is either cloning the full GATE repo and modifying existing components directly,
640 or, importing the `GATeNcOder` and `GATeMDeL` classes from GATE, and wrapping up
641 their model within it. Doing so requires the researcher to implement a relevant forward
642 function that can take in the modalities their model needs to process, as well as defining a
643 configuration that tells GATE what modalities a model can receive and output features on,
644 as well as any transforms needed for a batch to be ready for their model.
- 645 3. The user chooses a GATE tier to use (from `smallGATE`, `baseGATE` and `bigGATE`). Based
646 on the configuration defined by the user in step 2.
- 647 4. GATE generates a list of commands, each representing an experiment that needs to be run,
648 and can then run these commands on your local GPU box, parallelizing the tasks, one on
649 each available GPU, or, can provide a list of commands or json file that one can use to run
650 these commands on a GPU cluster, or other hardware.
- 651 5. GATE emits a wandb project, with metrics, visualizations and other measures, allowing easy
652 tracking of experiments, and sharing thereof, as well as huggingface model weights for each
653 model being trained – which is also used to achieve a *stateless* execution.
- 654 6. Once the experiments are completed, one can invoke the `produce-analysis.py` file within
655 GATE to get tables and figures that analyse the data, similar to what appears in this paper.
656 Those results can then be used to report results in a paper, or, be used to make decisions for
657 production models.

658 This process ensures the GATE framework is aware of what a model’s supported modalities are, as
659 well as how to produce modality-specific features, given the model. Once this is completed, the user,
660 with a single line of code, can select a GATE tier, and launch all jobs needed to produce results for that
661 tier. Importantly, GATE is made to facilitate and encourage foundation models that are diverse in their
662 capabilities, and allow the researchers to focus on what matters – that is, designing and training their
663 foundation model – rather than spending the majority of their time building and optimizing evaluation
664 boilerplate. Furthermore, the diversity of signal that GATE provides allows better understanding of a
665 given model’s strengths and weaknesses, which as a result makes the research, review and iteration
666 process of the field as a whole more efficient. This is because there is a consistent boilerplate that
667 runs all models, with broad signal that reduces probability of making erroneous conclusions – both in
668 the overly optimistic, or overly pessimistic side of things.

669 6.1 Principal Use Cases

- 670 1. **Model Development and Iteration:** GATE serves as a valuable tool during the model
671 research and development phase. By integrating the model into GATE and running either
672 the `smallGATE` or `baseGATE` tiers, developers can obtain a comprehensive and robust
673 performance evaluation of their model across diverse domains, tasks, and modalities. Worth
674 noting that GATE allows easy inclusion of foundation models **pretrained on images, video,
675 audio, text, etc.** to be **fine-tuned on pixel-based tasks**. It achieves this by replacing a
676 model’s root layer / embedding layer, with one appropriate for a given task’s modality, and
677 adding on top a relevant task adapter head.
- 678 2. **Model Evaluation for Machine Learning Research:** GATE enhances the communication
679 of research findings and their potential applications, a vital aspect of scientific collabora-
680 tion. By using GATE as a benchmark, even at the most cost-efficient GPU hour level
681 of `smallGATE`, the clarity and depth of future ML papers can be significantly improved.
682 GATE’s explicit evaluation of modality, domain, and task shifts in a given foundation model
683 provides a nuanced and informative perspective on a model’s true capabilities, offering a
684 more detailed understanding of a model’s strengths and weaknesses than optimizing a single
685 metric, such as ImageNet validation error.

686 7 Result Extras

687 The results were logged in WandB, and then further processed after all experiments were completed
688 to generate the tables and figures in this paper. Much of the logged information outside of testing
689 metrics were not used for any of the figures and tables in this paper. The full set of experiments and
690 all the logged results can be found at our wandb gate project repo².

691 7.1 Result Processing

692 Once all experiments were completed, we queried our wandb project repository and returned test
693 results from all our experiments, if an experiment name was duplicated, we used the latest entries,
694 and, for each experiment type there existed three independent runs. We averaged the results of any
695 metrics across such independent runs to acquire a better approximation to the true performance of
696 those models.

697 8 Preliminary Experiments Details

698 8.1 Preliminary Experiments

699 First, we trained models on ImageNet1k, CIFAR100, CLEVR, ADE20K, CityScapes, and, ACDC
700 for 5K iterations, using cosine annealing learning schedule or plateau annealing, with AdamW,
701 weight decays varying from 0.1 - 0.0001, and applied models from each major architecture category –
702 specifically, the CLIPViT, ImageNet pretrained ViT, ResNext, ResNet and ConvNextV2. The results
703 from these experiments pointed to the fact that there exists one general and good recipe for each
704 architecture style. The recipes that we discovered were as follows:

705 8.1.1 Across Architecture Settings

706 Unless otherwise stated, the settings here are applied universally in all experiments.

707 **Optimizer:** AdamW, weight decay 0.01, plateau annealing with patience 1000, relative scaling and
708 scale factor 0.5, and, threshold 0.0001.

709 **Training Details:** Training iterations: 10K, validate every 500 iterations.

710 **Test Details:** Top-3 validation models (across all validated checkpoints) are ensembled by prediction
711 averaging.

712 8.1.2 Architecture Specific Settings

713 **Convolutional Architectures: Optimizer:** AdamW, learning rate 1e-3, and for segmentation tasks
714 only, we used learning rate 6e-4

715 **Vision Transformer Architectures: Optimizer:** AdamW, learning rate 1e-5

716 **Convolutional + Transformer Hybrid Architectures Optimizer:** AdamW, learning rate 2e-5

717 The above recipes were what we used throughout all our experiments unless otherwise stated.

718 9 GATE Guiding Principles

719 The fundamental values driving the design decisions behind GATE are the following:

- 720 1. Maximizing Generalization Signal: GATE is designed to provide a high signal-to-noise
721 ratio concerning a model’s ability to generalize in diverse downstream contexts, that vary in
722 domain, task and modality. This allows for a more robust assessment of a model’s capacity
723 for adaptation and versatility. By noise here we refer to how clear a given signal response is.
724 For example, an image classification test accuracy signal on ImageNet, would provide clear

²omitted until double blind is over

725 signal with respect to the natural domain and the classification task, but would be blurry for
726 more compositional, object disentanglement and relational tasks, such as segmentation, or,
727 visual question answering.

728 2. Time Efficiency: Acknowledging the importance of computational resources and time,
729 GATE operates within set benchmarks of 12, 24, and 36 GPU hours (established on A100 @
730 40GB). These set timeframes ensure GATE’s assessments are both thorough and expedient.

731 3. Minimizing Usage Friction: The framework supporting GATE is designed to be user-friendly,
732 enabling easy integration of new backbones and facilitating smooth experimentation. This
733 low-friction approach ensures a streamlined experience when using GATE, making the
734 process of evaluation more efficient.

735 We argue that a good balance of the above can generate a pragmatic, yet thorough foundation model
736 evaluation suite, that will, importantly, be of real use to most researchers in the field.

737 10 Defining the GATE Benchmark

738 GATE is a comprehensive evaluation engine designed to advance the development of more general
739 machine learning models. It improves on existing benchmarks by enabling the evaluation of models
740 across diverse modalities, domains, and tasks.

741 GATE is composed of three key components. The first is a benchmark *pool*, a broad collection of
742 datasets, tasks, and processes that measure a model’s performance across various domains, tasks,
743 and modalities. The second component is a set of benchmark *tiers*, which are meticulously curated
744 subsets from the GATE benchmark pool, tailored to specific compute budgets and project phases.
745 The final, and is a software framework, designed to seamlessly integrate new foundation models and
746 execute the GATE tiers, thereby enabling efficient performance evaluation across a diverse range of
747 downstream modalities, domains, and tasks. Practically, GATE is directed towards machine learning
748 researchers and developers as a means to efficiently, and with little friction, get broad signal about
749 how their model performs after transfer in diverse contexts, specifically selected for their empirically
750 evaluated high signal-to-noise ratio with respect to predictive power in how a model performs in
751 previously unseen contexts.

752 Building GATE was a careful balancing act. We needed to respect specific time budgets while also
753 aiming for a wide variety of evaluation scenarios. Our approach was as follows:

- 754 1. Select a diverse set of learning contexts, spanning multiple domains, tasks and modalities.
755 We refer this as the *Benchmark Pool*.
- 756 2. Select a broad set of key foundation models, varying in their architecture, pretraining scheme
757 and source modality. We refer to this as the *Model Pool*.
- 758 3. Fine tune each of the models in the model pool, on each of the contexts in the benchmark
759 pool. Evaluate trained models on each context’s test sets.
- 760 4. Use the test set results acquired to quantify the predictive power each benchmark holds with
761 respect to previously unseen benchmarks, both at the individual level and the collection
762 level. We call this measure, the *downstream generalization predictability measure (DGPM)*.
- 763 5. Use the DGPM values of the various combinations of benchmarks to build the three GATE
764 tiers, selecting combinations of benchmarks that can provide the most information within a
765 target time budget.

766 We elaborate on each of the above steps in the following subsections.

767 11 Benchmark Pool Selection Details

768 **Medical Image Classification:** Medical data are known to present a substantial shift in both domain
769 and even modality depending on their format. We have selected datasets that not only pose significant
770 challenges for foundation models but also align with the broader imperative to deliver real-world
771 benefits downstream.

772 **Chexpert:** A dataset comprising a challenging array of chest x-rays annotated with findings critical to
773 diagnosing thoracic diseases. It tests models on their ability to navigate complex, multi-label medical
774 data, encapsulating the kind of nuanced decision-making that AI must augment in clinical settings.

775 **Diabetic Retinopathy Classification:** Early detection of diabetic retinopathy from retinal images
776 is a public health priority; models fine-tuned on this dataset can have immediate implications for
777 preventing vision loss on a global scale. This dataset requires models to decipher fine-grained,
778 progressive changes indicative of the disease, reflecting the precision necessary for medical AI
779 applications.

780 **HAM10000 (Human Against Machine with 10000 dermoscopic images):** The dataset provides
781 a diverse spectrum of skin lesion images vital for differentiating between benign and malignant
782 conditions. Incorporating this dataset not only challenges the pattern recognition prowess of AI but
783 also contributes to the advancement of dermatology through machine learning technologies.

784 **Metrics:** We collect Average Precision Score (APS), Area Under the Receiver Operating Char-
785 acteristics Curve (AUC), and Brier Score (BS) both overall (i.e. macro) as well as for individual
786 pathologies/classes.

787 **Medical Segmentation:** This category evaluates foundational models’ ability to generalize from
788 natural to medical image modalities and to perform domain-specific tasks that require precision and
789 complex spatial understanding:

790 **ACDC (Automated Cardiac Diagnosis Challenge):** This dataset is aimed at assessing models’
791 generalization to the medical domain, particularly the transferability of representations for segmenting
792 anatomical structures in cardiac MRI images. By focusing on the heart’s intricate anatomy, ACDC
793 tests the models’ ability to adapt to clinically relevant shapes and patterns—a shift from common
794 visual recognition tasks to precise medical delineation. **Metrics:** We collect dice loss, mIoU, mean
795 accuracy and overall accuracy.

796 12 Benchmark Pool Details

797 Having a set of diverse benchmarks ranging in challenge factor, as well as modality, task and domain
798 shift was key. We explain in more detail why we consider these factors important in Appendix in
799 more detail. We refer to this as our *benchmark pool*, and it consists of the following:

800 **Image Classification:** We employ **ImageNet1k** [9], **CIFAR100** [28], **Places365** [74], and **Food101**
801 [36] to cover diverse natural image domains. Additionally, we include **HappyWhale** [17] for a more
802 challenging domain shift, aiding in wildlife research and providing an interesting test case for model
803 evaluation.

804 **Few Shot Image Classification:** We use the MetaDataset task recipe on the **Aircraft** [37], **Fungi**
805 [50], **MiniImageNet** [62], **CUB200** [63], and **Describable Features** [69] datasets to evaluate task
806 and domain shift robustness for an evaluation model.

807 **Zero Shot Text-Image Classification:** Another key setting is that of zero-shot text-image classifica-
808 tion, on which many current key models were trained and evaluated [43]. We utilize **Flickr30K**, **New**
809 **Yorker Caption Context** (a challenging humor task), and **Winoground**—a task requiring the model
810 to match two texts with their corresponding images, focusing on compositional differences.

811 **Visual Relational Reasoning:** A context where earlier models, such as ResNet50 [18] had low
812 performance without layers with associative inductive biases (e.g., relational neural networks or
813 transformers [49, 61]). This ensures we are aware of any trade-offs in relational compositional
814 abilities in our models. We use **CLEVR** [23] and **CLEVRMath** [34].

815 **Image Semantic Segmentation:** Essential for various real-world applications, serving as an indicator
816 of a model’s ability to retain spatial information and identify objects at a per-pixel level. **ADE20K**
817 [75], **COCO10K** [33], **COCO164K** [33], **NYU-Depth-v2** [54], **PascalContext** [38], and **Cityscapes**
818 [8].

819 **Medical Image Classification:** Medical data exhibit substantial domain and modality shifts, posing
820 significant challenges for machine learning models while aligning with the imperative to deliver
821 real-world benefits. **Chexpert** [21] (chest X-rays annotated for thoracic disease diagnosis), **Dia-**

822 **abetic Retinopathy Classification** [16] (retinal images for early detection of diabetic retinopathy),
823 **HAM10000** [60] (dermatoscopic images for differentiating skin lesions).

824 **Medical Segmentation** → **ACDC (Automated Cardiac Diagnosis Challenge)** [5]: This dataset as-
825 sesses models’ generalization to the medical domain, particularly the transferability of representations
826 for segmenting anatomical structures in cardiac MRI images. By focusing on the heart’s intricate
827 anatomy, ACDC tests the models’ ability to adapt to clinically relevant shapes and patterns.

828 **Video Classification:** Video classification tasks test models on their temporal generalization abilities
829 and require an understanding of not only individual frame content but also the transition and context
830 between frames. **HMDB51 (Human Motion Database)** [30], **UCF-101 (University of Central**
831 **Florida - 101 action categories)** [56], **Kinetics400** [24].

832 **Video Regression:** Where classification tasks gauge categorical distinctions, video regression tasks
833 assess models’ ability to make continuous numerical predictions from temporal data, serving as an
834 indicator of a model’s capability to process and quantify dynamic content. **iWildcam (International**
835 **Wildlife Camera Trap Challenge)** [4]: This dataset targets estimating animal species abundance from
836 videos and is a direct test of modality and task shift, and showcases a models’ potential impact on
837 ecological monitoring and species conservation efforts.

- 838 1. **Modality shifting** contexts: Contexts where the foundation model is asked to learn to do
839 well at a task that requires understanding of a previously unseen modality. More specifically,
840 assuming a foundation model has been trained on natural images, this would be transferring
841 to medical imaging, video, audio and test contexts. This would shed light on the performance
842 of a model’s middle layers.
- 843 2. **Task shifting** contexts: Contexts where a model is tasked with performing a previously
844 unseen task, for example, transferring from classification to segmentation or relational
845 reasoning.
- 846 3. **Domain shifting** contexts: Contexts where a model is required to perform a task on a
847 domain that is different from the one it was trained on. For example moving from natural
848 images on ImageNet at 224x224 resolution to black and white Omniglot characters at 28x28
849 resolution, or, moving from ImageNet to images of fungi. More extreme domain shifts
850 would be going from natural images to medical images for example.

851 **13 Model Pool Details**

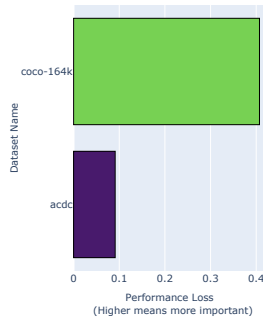
852 **14 Task Adapter Details**

853 **15 Experimental Details**

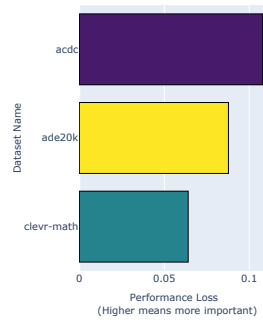
854 **Experimental Environment Details:** GPUs: 4 x A6000 Ada @ 48GB, CPUs: 128 Core AMD
855 EPYC 7713 64-Core Processor, RAM: 1 TB, HD: 15TB NVME. All experiments were done with
856 BF16 precision.

857 **16 Additional Results**

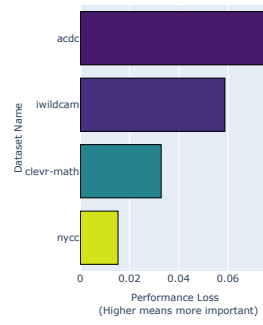
858 **16.1 Full details on discovered combinations**



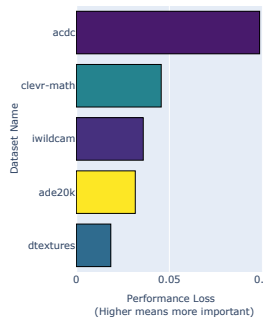
(a) Best $k=2$ discovered combination



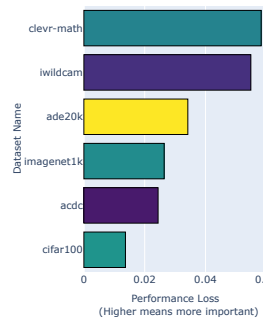
(b) Best $k=3$ discovered combination



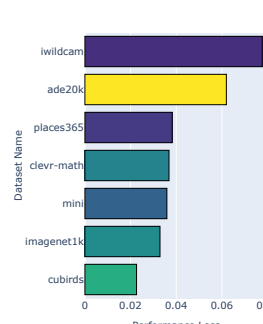
(c) Best $k=4$ discovered combination



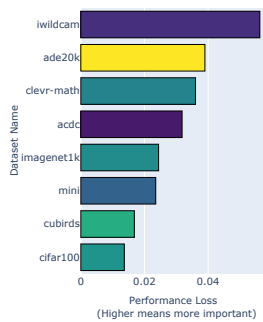
(d) Best $k=5$ discovered combination



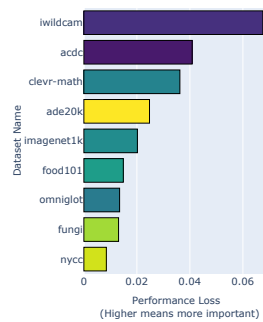
(e) Best $k=6$ discovered combination



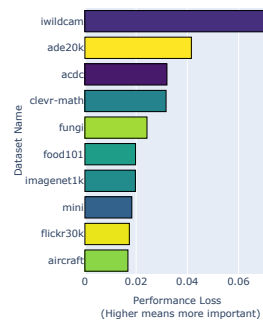
(f) Best $k=7$ discovered combination



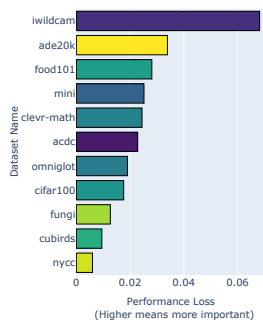
(g) Best $k=8$ discovered combination



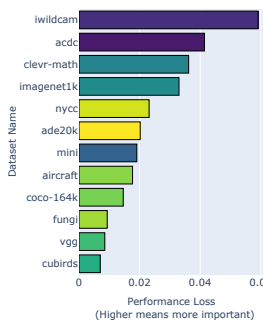
(h) Best $k=9$ discovered combination



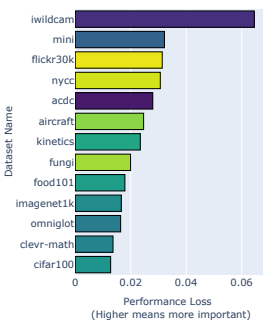
(i) Best $k=10$ discovered combination



(j) Best $k=11$ discovered combination



(k) Best $k=12$ discovered combination

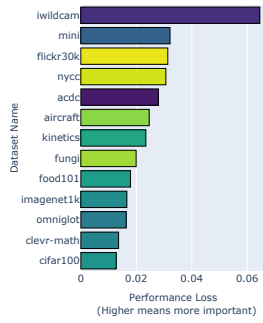


(l) Best $k=13$ discovered combination

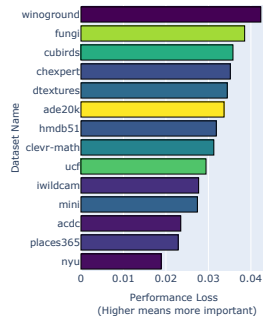
Figure 6: Degradation of predictive power when a given benchmark is removed and the meta-model trained from scratch, for different best combinations in varying k .

Kinetics Acc@1	48.8	44.2	51.4	43.7	40.3	44.6	33.2	36.4	25.8	2.7	1.0	0.2	0.3	0.4	2.0	1.6	1.0	0.5	0.3	0.3	0.3
Kinetics Acc@5	75.5	70.9	77.9	70.7	67.6	71.7	59.9	63.0	51.8	9.7	4.3	1.3	1.4	1.7	7.0	6.5	3.5	2.2	1.3	1.3	1.3
Kinetics Loss	2.4	2.6	2.1	2.5	2.7	2.5	3.2	3.0	3.5	5.5	6.1	6.1	6.1	6.1	5.7	5.8	6.0	6.1	6.1	6.1	6.1
UCF-101 Acc@1	84.4	75.1	69.9	63.2	75.0	63.4	58.8	66.6	48.7	19.7	11.1	2.8	0.8	2.1	15.2	13.3	6.6	8.7	6.5	7.0	2.7
UCF-101 Acc@5	95.4	92.5	89.1	82.3	91.6	86.2	81.7	86.3	75.3	42.2	28.9	8.5	5.0	8.2	35.5	33.8	17.9	25.2	23.1	20.2	11.2
UCF-101 Loss	0.6	1.0	1.3	1.7	1.0	1.5	1.7	1.4	2.3	4.3	5.0	4.8	4.7	4.6	3.7	3.8	4.5	4.0	4.2	4.2	4.5
Task Mean	73.0	65.6	66.6	58.8	63.7	57.5	53.3	57.5	49.8	17.2	14.3	4.2	3.3	5.0	15.7	14.7	8.8	10.9	10.3	10.2	5.8
Video Reg																					
IWildCam MAE Score	1.3	1.4	1.3	1.4	1.4	1.6	1.4	1.5	1.6	2.0	1.9	1.9	2.6	2.1	1.8	1.8	1.9	1.8	2.2	1.8	2.1
IWildCam MSE Loss	3.7	4.4	4.0	4.0	4.1	5.4	4.3	5.0	5.9	7.1	6.5	6.2	12.5	8.5	5.1	6.3	6.0	6.2	8.6	6.4	8.4
Task Mean	1.3	1.4	1.3	1.4	1.4	1.6	1.4	1.5	1.6	2.0	1.9	1.9	2.6	2.1	1.8	1.8	1.9	1.8	2.2	1.8	2.1
GATE																					
Full GATE Mean	69.0	66.8	66.8	64.6	64.3	63.4	62.1	62.2	58.5	56.3	54.4	48.4	42.8	39.6	37.5	37.2	36.2	36.9	35.0	34.9	31.8
Big GATE Mean	76.6	74.5	74.4	72.8	72.0	71.9	70.6	70.0	66.8	66.7	64.8	58.5	53.1	46.8	43.8	43.4	41.9	41.5	40.9	39.8	37.1
Base GATE Mean	68.3	65.6	65.7	62.6	63.7	60.7	60.2	60.7	58.6	55.1	53.5	48.2	42.8	38.0	36.5	36.3	35.4	36.6	34.8	34.8	30.4
Small GATE Mean	77.7	74.9	74.6	73.3	72.4	71.2	68.9	69.1	65.3	65.7	61.7	58.5	49.3	40.5	35.7	35.4	35.9	35.3	34.1	34.4	30.4
Full GATE Rank	1.0	3.0	2.0	4.0	5.0	6.0	8.0	7.0	9.0	10.0	11.0	12.0	13.0	14.0	15.0	16.0	18.0	17.0	19.0	20.0	21.0
Big GATE Rank	1.0	2.0	3.0	4.0	5.0	6.0	7.0	8.0	9.0	10.0	11.0	12.0	13.0	14.0	15.0	16.0	17.0	18.0	19.0	20.0	21.0
Base GATE Rank	1.0	3.0	2.0	5.0	4.0	7.0	8.0	6.0	9.0	10.0	11.0	12.0	13.0	14.0	16.0	17.0	18.0	15.0	20.0	19.0	21.0
Small GATE Rank	1.0	2.0	3.0	4.0	5.0	6.0	8.0	7.0	10.0	9.0	11.0	12.0	13.0	14.0	16.0	17.0	15.0	18.0	20.0	19.0	21.0

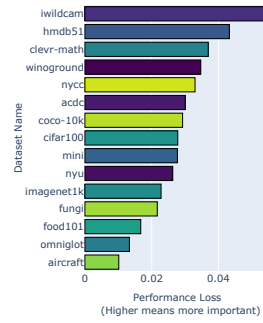
Table 3: Full experiments table: Black/Bold best model, Green second best, Blue third best, and red the worst performing model. Models prefixed with 's' refer to 'from scratch' trained models, rather than pretrained. This table showcases the full set of data we use to evolve GATE using EEVEE.



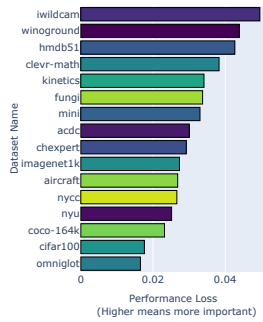
(a) Best $k=13$ discovered combination



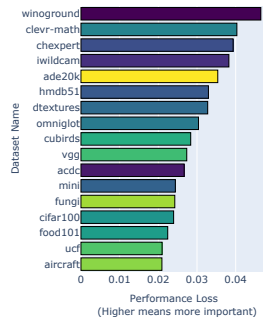
(b) Best $k=14$ discovered combination



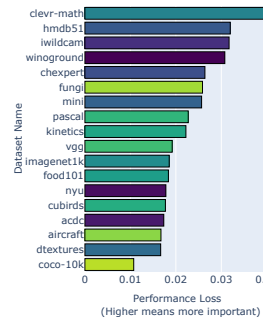
(c) Best $k=15$ discovered combination



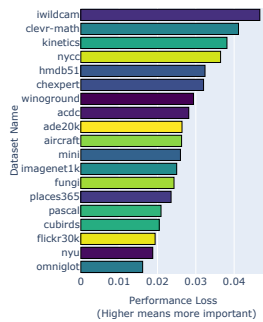
(d) Best $k=16$ discovered combination



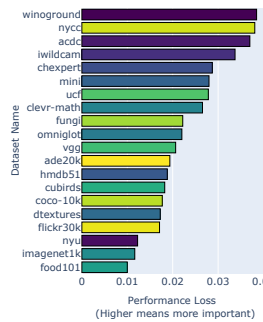
(e) Best $k=17$ discovered combination



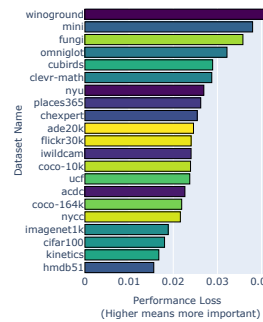
(f) Best $k=18$ discovered combination



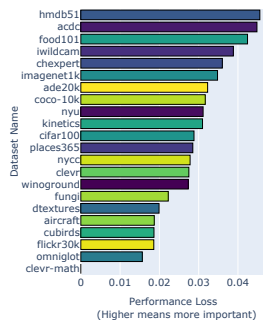
(g) Best $k=19$ discovered combination



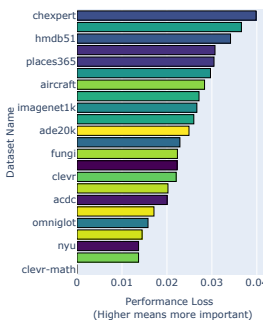
(h) Best $k=20$ discovered combination



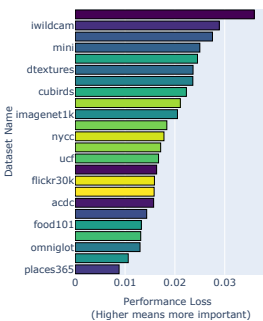
(i) Best $k=21$ discovered combination



(j) Best $k=22$ discovered combination

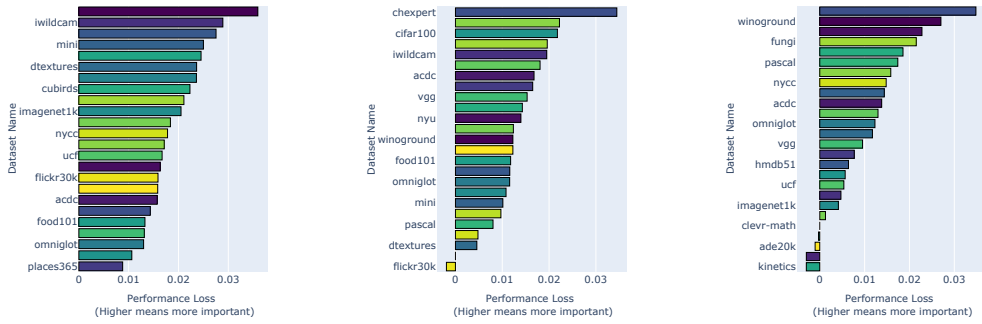


(k) Best $k=23$ discovered combination

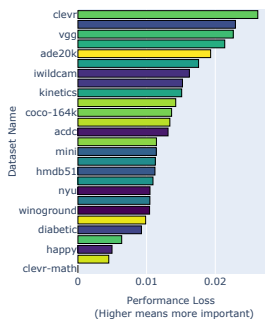


(l) Best $k=24$ discovered combination

Figure 7: Degradation of predictive power when a given benchmark is removed and the meta-model trained from scratch, for different best combinations in varying k .



(a) Best $k=24$ discovered combination (b) Best $k=25$ discovered combination (c) Best $k=26$ discovered combination



(d) Best $k=27$ discovered combination

Figure 8: Degradation of predictive power when a given benchmark is removed and the meta-model trained from scratch, for different best combinations in varying k .

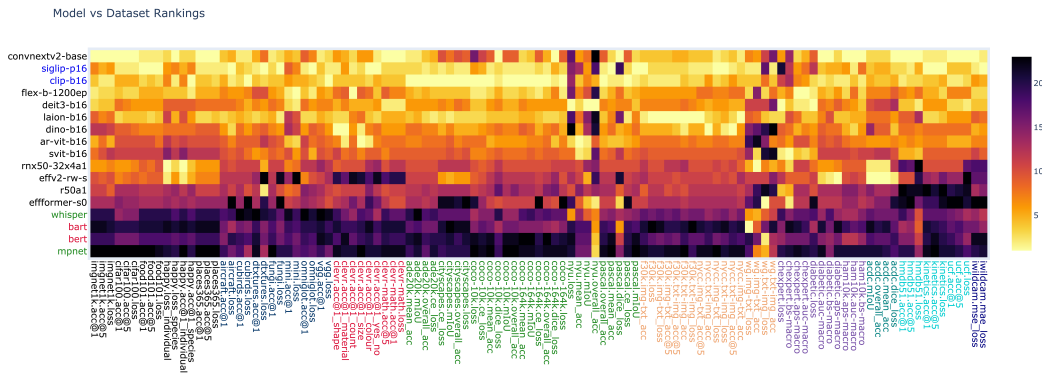


Figure 9: Ranking Heatmap for bigGATE We show how the various models on the y-axis rank on the metrics on the x-axis, where brighter is higher/better rank. From left to right we apply a spearman correlation sorting to capture tasks more similar to imagenet1k more towards the leftmost side, and, dissimilar ones towards the rightmost side. From top to bottom we rank models based on average rank.

Model vs Dataset Rankings

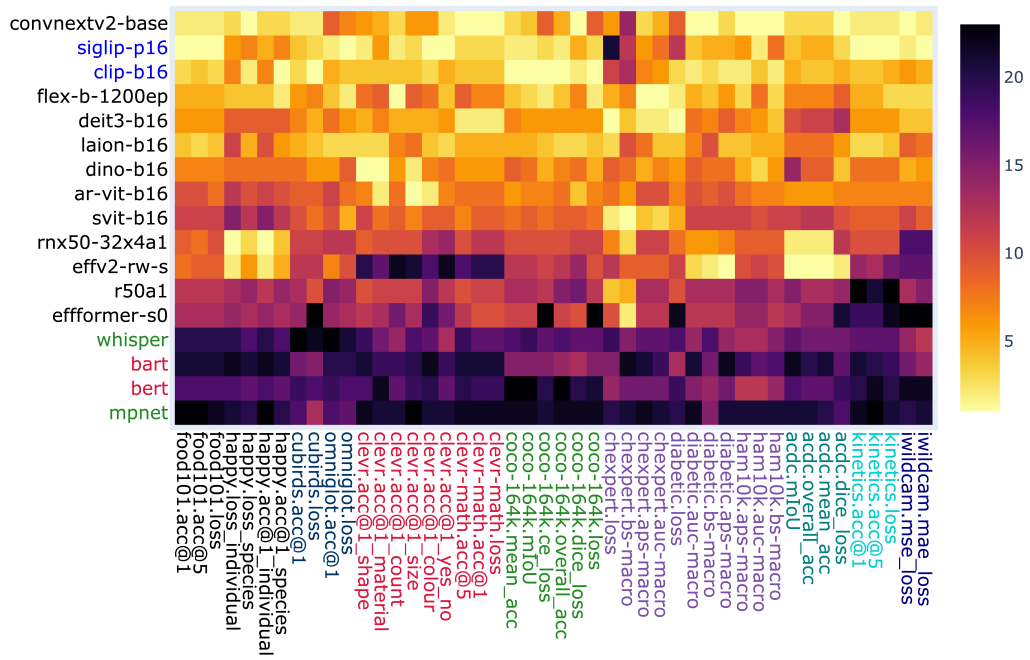


Figure 10: Ranking Heatmap for baseGATE: We show how the various models on the y-axis rank on the metrics on the x-axis, where brighter is higher/better rank. From left to right we apply a spearman correlation sorting to capture tasks more similar to imagenet1k more towards the leftmost side, and, dissimilar ones towards the rightmost side. From top to bottom we rank models based on average rank.

Model vs Dataset Rankings

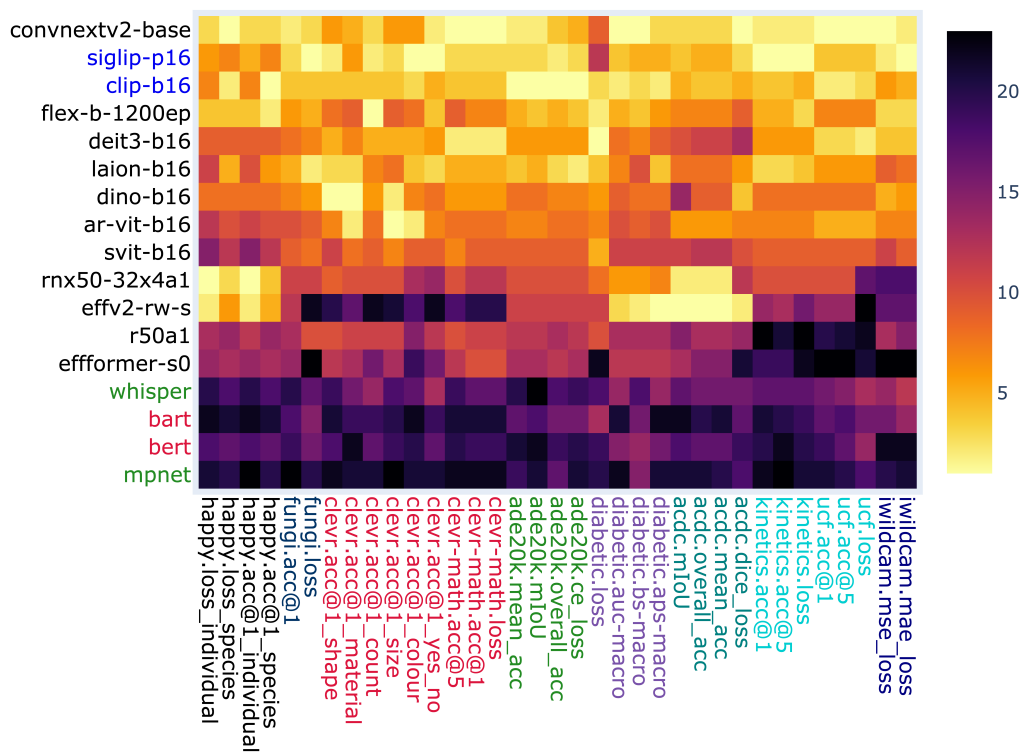


Figure 11: Ranking Heatmap for smallGATE: We show how the various models on the y-axis rank on the metrics on the x-axis, where brighter is higher/better rank. From left to right we apply a spearman correlation sorting to capture tasks more similar to imagenet1k more towards the leftmost side, and, dissimilar ones towards the rightmost side. From top to bottom we rank models based on average rank.

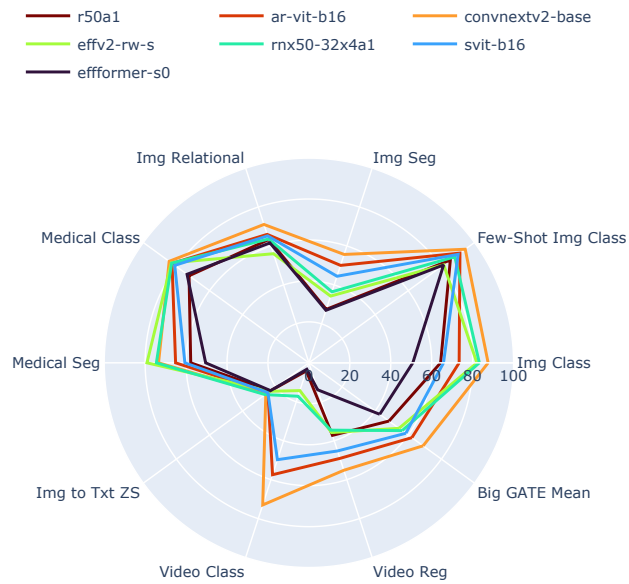


Figure 12: Architecture Variation: Results of keeping the pretraining method the same as ImageNet1k classification and varying the architecture across various key task domains.

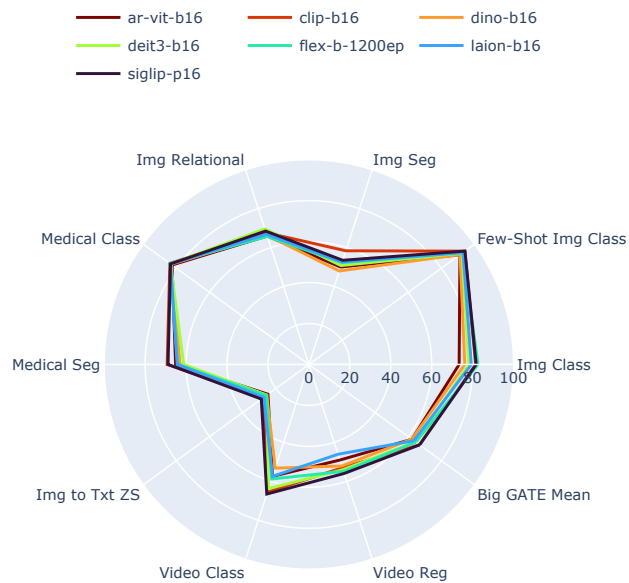


Figure 13: Pretraining Scheme Variation: Results of varying the pretraining method and keeping the architecture as ViT B16 across various key task domains.



Figure 14: Modality Variation: Results of attempting modality shifting from audio and text to vision tasks.

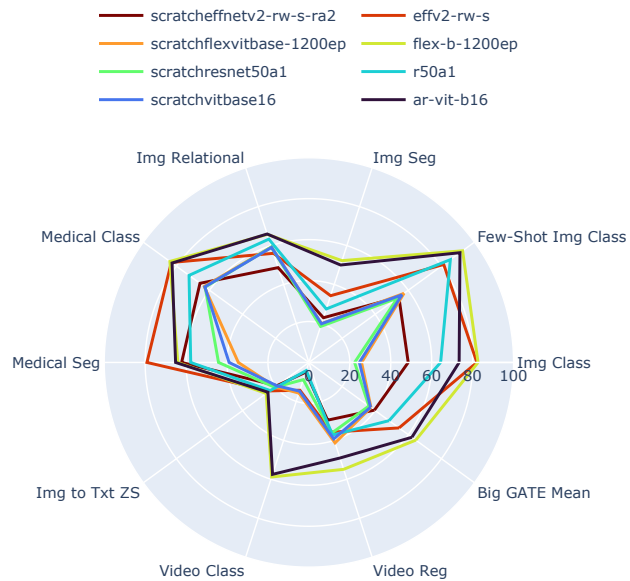


Figure 15: Modality Variation: Results of attempting modality shifting from audio and text to vision tasks.

859 **NeurIPS Paper Checklist**

860 **1. Claims**

861 Question: Do the main claims made in the abstract and introduction accurately reflect the
862 paper's contributions and scope?

863 Answer: [Yes]

864 Justification: All the claims made are substantiated with rigorous empirical results and
865 communicated via tables and figures.

866 Guidelines:

- 867 • The answer NA means that the abstract and introduction do not include the claims
868 made in the paper.
- 869 • The abstract and/or introduction should clearly state the claims made, including the
870 contributions made in the paper and important assumptions and limitations. A No or
871 NA answer to this question will not be perceived well by the reviewers.
- 872 • The claims made should match theoretical and experimental results, and reflect how
873 much the results can be expected to generalize to other settings.
- 874 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
875 are not attained by the paper.

876 **2. Limitations**

877 Question: Does the paper discuss the limitations of the work performed by the authors?

878 Answer: [Yes]

879 Justification: We have an explicit limitations section.

880 Guidelines:

- 881 • The answer NA means that the paper has no limitation while the answer No means that
882 the paper has limitations, but those are not discussed in the paper.
- 883 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 884 • The paper should point out any strong assumptions and how robust the results are to
885 violations of these assumptions (e.g., independence assumptions, noiseless settings,
886 model well-specification, asymptotic approximations only holding locally). The authors
887 should reflect on how these assumptions might be violated in practice and what the
888 implications would be.
- 889 • The authors should reflect on the scope of the claims made, e.g., if the approach was
890 only tested on a few datasets or with a few runs. In general, empirical results often
891 depend on implicit assumptions, which should be articulated.
- 892 • The authors should reflect on the factors that influence the performance of the approach.
893 For example, a facial recognition algorithm may perform poorly when image resolution
894 is low or images are taken in low lighting. Or a speech-to-text system might not be
895 used reliably to provide closed captions for online lectures because it fails to handle
896 technical jargon.
- 897 • The authors should discuss the computational efficiency of the proposed algorithms
898 and how they scale with dataset size.
- 899 • If applicable, the authors should discuss possible limitations of their approach to
900 address problems of privacy and fairness.
- 901 • While the authors might fear that complete honesty about limitations might be used by
902 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
903 limitations that aren't acknowledged in the paper. The authors should use their best
904 judgment and recognize that individual actions in favor of transparency play an impor-
905 tant role in developing norms that preserve the integrity of the community. Reviewers
906 will be specifically instructed to not penalize honesty concerning limitations.

907 **3. Theory Assumptions and Proofs**

908 Question: For each theoretical result, does the paper provide the full set of assumptions and
909 a complete (and correct) proof?

910 Answer: [NA]

911 Justification: No theories were derived.

912 Guidelines:

- 913 • The answer NA means that the paper does not include theoretical results.
- 914 • All the theorems, formulas, and proofs in the paper should be numbered and cross-
- 915 referenced.
- 916 • All assumptions should be clearly stated or referenced in the statement of any theorems.
- 917 • The proofs can either appear in the main paper or the supplemental material, but if
- 918 they appear in the supplemental material, the authors are encouraged to provide a short
- 919 proof sketch to provide intuition.
- 920 • Inversely, any informal proof provided in the core of the paper should be complemented
- 921 by formal proofs provided in appendix or supplemental material.
- 922 • Theorems and Lemmas that the proof relies upon should be properly referenced.

923 4. Experimental Result Reproducibility

924 Question: Does the paper fully disclose all the information needed to reproduce the main ex-

925 perimental results of the paper to the extent that it affects the main claims and/or conclusions

926 of the paper (regardless of whether the code and data are provided or not)?

927 Answer: [Yes]

928 Justification: We do so both in the main paper, and in more detail in the appendix, in addition

929 to offering the codebase that reproduces all results.

930 Guidelines:

- 931 • The answer NA means that the paper does not include experiments.
- 932 • If the paper includes experiments, a No answer to this question will not be perceived
- 933 well by the reviewers: Making the paper reproducible is important, regardless of
- 934 whether the code and data are provided or not.
- 935 • If the contribution is a dataset and/or model, the authors should describe the steps taken
- 936 to make their results reproducible or verifiable.
- 937 • Depending on the contribution, reproducibility can be accomplished in various ways.
- 938 For example, if the contribution is a novel architecture, describing the architecture fully
- 939 might suffice, or if the contribution is a specific model and empirical evaluation, it may
- 940 be necessary to either make it possible for others to replicate the model with the same
- 941 dataset, or provide access to the model. In general, releasing code and data is often
- 942 one good way to accomplish this, but reproducibility can also be provided via detailed
- 943 instructions for how to replicate the results, access to a hosted model (e.g., in the case
- 944 of a large language model), releasing of a model checkpoint, or other means that are
- 945 appropriate to the research performed.
- 946 • While NeurIPS does not require releasing code, the conference does require all submis-
- 947 sions to provide some reasonable avenue for reproducibility, which may depend on the
- 948 nature of the contribution. For example
 - 949 (a) If the contribution is primarily a new algorithm, the paper should make it clear how
 - 950 to reproduce that algorithm.
 - 951 (b) If the contribution is primarily a new model architecture, the paper should describe
 - 952 the architecture clearly and fully.
 - 953 (c) If the contribution is a new model (e.g., a large language model), then there should
 - 954 either be a way to access this model for reproducing the results or a way to reproduce
 - 955 the model (e.g., with an open-source dataset or instructions for how to construct
 - 956 the dataset).
 - 957 (d) We recognize that reproducibility may be tricky in some cases, in which case
 - 958 authors are welcome to describe the particular way they provide for reproducibility.
 - 959 In the case of closed-source models, it may be that access to the model is limited in
 - 960 some way (e.g., to registered users), but it should be possible for other researchers
 - 961 to have some path to reproducing or verifying the results.

962 5. Open access to data and code

963 Question: Does the paper provide open access to the data and code, with sufficient instruc-

964 tions to faithfully reproduce the main experimental results, as described in supplemental

965 material?

966
967
968
969
970
971
972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1010
1011
1012
1013
1014
1015
1016

Answer: [Yes]

Justification: Full code and data are available and shared on github and huggingface.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

6. Experimental Setting/Details

Question: Does the paper specify all the training and test details (e.g., data splits, hyper-parameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: We describe these in the experiments section in summary, and in the appendix in detail.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

7. Experiment Statistical Significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [Yes]

Justification: Where relevant our results include error bars.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).

- 1017 • It should be clear whether the error bar is the standard deviation or the standard error
1018 of the mean.
- 1019 • It is OK to report 1-sigma error bars, but one should state it. The authors should
1020 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis
1021 of Normality of errors is not verified.
- 1022 • For asymmetric distributions, the authors should be careful not to show in tables or
1023 figures symmetric error bars that would yield results that are out of range (e.g. negative
1024 error rates).
- 1025 • If error bars are reported in tables or plots, The authors should explain in the text how
1026 they were calculated and reference the corresponding figures or tables in the text.

1027 8. Experiments Compute Resources

1028 Question: For each experiment, does the paper provide sufficient information on the com-
1029 puter resources (type of compute workers, memory, time of execution) needed to reproduce
1030 the experiments?

1031 Answer: [TODO]

1032 Justification: [TODO]

1033 Guidelines:

- 1034 • The answer NA means that the paper does not include experiments.
- 1035 • The paper should indicate the type of compute workers CPU or GPU, internal cluster,
1036 or cloud provider, including relevant memory and storage.
- 1037 • The paper should provide the amount of compute required for each of the individual
1038 experimental runs as well as estimate the total compute.
- 1039 • The paper should disclose whether the full research project required more compute
1040 than the experiments reported in the paper (e.g., preliminary or failed experiments that
1041 didn't make it into the paper).

1042 9. Code Of Ethics

1043 Question: Does the research conducted in the paper conform, in every respect, with the
1044 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

1045 Answer: [Yes]

1046 Justification: Yes it does abide by the code of ethics to our best of our understanding.

1047 Guidelines:

- 1048 • The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- 1049 • If the authors answer No, they should explain the special circumstances that require a
1050 deviation from the Code of Ethics.
- 1051 • The authors should make sure to preserve anonymity (e.g., if there is a special consid-
1052 eration due to laws or regulations in their jurisdiction).

1053 10. Broader Impacts

1054 Question: Does the paper discuss both potential positive societal impacts and negative
1055 societal impacts of the work performed?

1056 Answer: [No]

1057 Justification: It's a method for finding optimal subsets of benchmarks from a large pool and
1058 a framework that automates model encoder evaluation. Societal impacts relate to improved
1059 research efficiency and hopefully compute usage, however this is too far from what one
1060 would consider strongly tied societal impacts.

1061 Guidelines:

- 1062 • The answer NA means that there is no societal impact of the work performed.
- 1063 • If the authors answer NA or No, they should explain why their work has no societal
1064 impact or why the paper does not address societal impact.
- 1065 • Examples of negative societal impacts include potential malicious or unintended uses
1066 (e.g., disinformation, generating fake profiles, surveillance), fairness considerations
1067 (e.g., deployment of technologies that could make decisions that unfairly impact specific
1068 groups), privacy considerations, and security considerations.

- 1069
- 1070
- 1071
- 1072
- 1073
- 1074
- 1075
- 1076
- 1077
- 1078
- 1079
- 1080
- 1081
- 1082
- 1083
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
 - The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
 - If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

1084 11. Safeguards

1085 Question: Does the paper describe safeguards that have been put in place for responsible
1086 release of data or models that have a high risk for misuse (e.g., pretrained language models,
1087 image generators, or scraped datasets)?

1088 Answer: [NA]

1089 Justification: It's a benchmark with datasets that are already public and previously published
1090 in other papers.

1091 Guidelines:

- 1092
- 1093
- 1094
- 1095
- 1096
- 1097
- 1098
- 1099
- 1100
- 1101
- The answer NA means that the paper poses no such risks.
 - Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
 - Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
 - We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

1102 12. Licenses for existing assets

1103 Question: Are the creators or original owners of assets (e.g., code, data, models), used in
1104 the paper, properly credited and are the license and terms of use explicitly mentioned and
1105 properly respected?

1106 Answer: [Yes]

1107 Justification: All datasets used have appropriate licenses, and the code packages used in
1108 implementing our software framework have appropriate licenses as well.

1109 Guidelines:

- 1110
- 1111
- 1112
- 1113
- 1114
- 1115
- 1116
- 1117
- 1118
- 1119
- 1120
- The answer NA means that the paper does not use existing assets.
 - The authors should cite the original paper that produced the code package or dataset.
 - The authors should state which version of the asset is used and, if possible, include a URL.
 - The name of the license (e.g., CC-BY 4.0) should be included for each asset.
 - For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
 - If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.

- 1121 • For existing datasets that are re-packaged, both the original license and the license of
1122 the derived asset (if it has changed) should be provided.
1123 • If this information is not available online, the authors are encouraged to reach out to
1124 the asset’s creators.

1125 13. New Assets

1126 Question: Are new assets introduced in the paper well documented and is the documentation
1127 provided alongside the assets?

1128 Answer: [Yes]

1129 Justification: Our codebase is fully documented.

1130 Guidelines:

- 1131 • The answer NA means that the paper does not release new assets.
1132 • Researchers should communicate the details of the dataset/code/model as part of their
1133 submissions via structured templates. This includes details about training, license,
1134 limitations, etc.
1135 • The paper should discuss whether and how consent was obtained from people whose
1136 asset is used.
1137 • At submission time, remember to anonymize your assets (if applicable). You can either
1138 create an anonymized URL or include an anonymized zip file.

1139 14. Crowdsourcing and Research with Human Subjects

1140 Question: For crowdsourcing experiments and research with human subjects, does the paper
1141 include the full text of instructions given to participants and screenshots, if applicable, as
1142 well as details about compensation (if any)?

1143 Answer: [NA]

1144 Justification: No crowdsourcing with humans

1145 Guidelines:

- 1146 • The answer NA means that the paper does not involve crowdsourcing nor research with
1147 human subjects.
1148 • Including this information in the supplemental material is fine, but if the main contribu-
1149 tion of the paper involves human subjects, then as much detail as possible should be
1150 included in the main paper.
1151 • According to the NeurIPS Code of Ethics, workers involved in data collection, curation,
1152 or other labor should be paid at least the minimum wage in the country of the data
1153 collector.

1154 15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human 1155 Subjects

1156 Question: Does the paper describe potential risks incurred by study participants, whether
1157 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)
1158 approvals (or an equivalent approval/review based on the requirements of your country or
1159 institution) were obtained?

1160 Answer: [NA]

1161 Justification: Same as previous answer.

1162 Guidelines:

- 1163 • The answer NA means that the paper does not involve crowdsourcing nor research with
1164 human subjects.
1165 • Depending on the country in which research is conducted, IRB approval (or equivalent)
1166 may be required for any human subjects research. If you obtained IRB approval, you
1167 should clearly state this in the paper.
1168 • We recognize that the procedures for this may vary significantly between institutions
1169 and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the
1170 guidelines for their institution.
1171 • For initial submissions, do not include any information that would break anonymity (if
1172 applicable), such as the institution conducting the review.