Side-Tuning: Network Adaptation via Additive Side Networks

Jeffrey O. Zhang1* Alexander Sax1* Amir Zamir1,2

University of California, Berkeley 2Stanford University 3Facebook AI Research

Fixed Features | Fine-Tune | Side-Tune

that integrates the B(x) prior with the evidence from S(x).

When training a neural network for a desired task, one may prefer to adapt a pretrained network rather than a randomly initialized one - due to lacking enough training data, performing lifelong learning where the system has to learn a new task while being previously trained for other tasks, or wishing to encode priors in the network via preset Pigure 1. The side-tuning framework vs the common alternawork adaptation are fine-tuning and using the pre-trained work that should be adapted to a new task, fine-tuning re-trains the inter-task interference is not possible. network as a fixed feature extractor, among others. pretrained network's weights and fixed feature extraction trains a

In this paper we propose a straightforward alternative: readout function with no re-training of the pre-trained weights. In Side-Tuning, Side-tuning adapts a pre-trained network by contrast, Side-tuning adapts the pre-trained network by training a training a lightweight "side" network that is fused with lightweight conditioned "side" network that is fused with the (unthe (unchanged) pre-trained network using a simple additive process. This simple method works as well as or better than existing solutions while it resolves some of the a side model $S: \mathbb{X} \to \mathbb{Y}$, so that the representations for the basic issues with fine-tuning, fixed features, and several target task are other common baselines. In particular, side-tuning is less prone to overfitting when little training data is available. vields better results than using a fixed feature extractor. and doesn't suffer from catastrophic forgetting in lifelong for some combining operation . We use a learned alphalearning. We demonstrate the performance of side-tuning blending, $a \oplus b \triangleq \alpha a + (1 - \alpha)b$ for this purpose (other under a diverse set of scenarios, including lifelong learn- options are discussed in Section 3.1). ing (iCIFAR, Taskonomy), reinforcement learning, imitation learning (visual navigation in Habitat). NLP question- framework to: fine-tuning, feature extration, and stage-wise answering (SQuAD v2), and single-task transfer learning training (see Fig. 2, right). Hence those can be viewed as (Taskonomy) with consistently promising results special cases of the general side-tuning framework. Other curricula suggest (e.g.) a maximum a posteriori estimator

Side-tuning is an example of an additive learning approach as it adds (strategically placed) parameters for each The goal of side-tuning is to capitalize on a pretrained new task. Fixed feature extraction would be a simple exmodel to better learn one or more novel tasks. By design, ample of an additive approach with zero new parameters. side-tuning does so without degrading performance of the As a result, fixed features are don't adapt the base network base model. The framework is straightforward: it assumes over the lifetime of the agent. A number of existing apaccess to the frozen base model $B: X \to Y$ that maps inproaches address this by learning new parameters (the numputs into some representation space that is shared between ber of which scales with the size of the base network) for in \mathbb{R}^N) or actual model predictions. Side-tuning then learns

leads to overfitting. Side-tuning is a simple method that addresses these limitations by strategically placing a small number of learnable parameters.

base requires minor updates. By adding fewer parameters Incremental learning's objective is to learn a sequence per task, side-tuning can learn more tasks before the model of tasks T. T., and perform well on the entire set at the grows large enough to require parameter consolidation. end of training. Two problems arise from this sequential These approaches stand in contrast to most existing presentation: catastrophic forgetting (see Sec. 4.2.1) and

methods for incremental learning, which do not increase the learning speed (see Sec. 4.2.2). Incremental learning meth-

existing methods constrain the learning procedure, leading nod Low Data High Data (Lifelong) to undesirable trade-offs. By design, additive approaches are able to circumvent forgetting, reuse knowledge, and scale to more tasks. Side-tuning is one of the simplest con-Table 1. Advantages of side-tuning vs. alternatives. Fixed features has sistent additive methods which takes this family of methno learnable parameters and thus cannot adapt to new information. On the ods and captures a small yet core component that makes other hand, fine-tuning has a large number of learnable parameters which them work. We show this experimentally on various tasks v2 [25], and Taskonomy [36]. In the remainder of this section we overview sidetuning's connection to related fields.

and datasets, including iCIFAR [27], Habitat [32], SOuAD

than the alternative methods.

- nections from neural networks of earlier tasks. Side-tuning learning new information. Discussed in Sec. 4.2.1.

highly performant, which we demonstrate in Section 4.2.3. Moreover, it is significantly simpler than most existing life-

This simple approach deals with the key challenges of

Broadly speaking, network adaptation methods either the motivation for our method.

number of parameters over time and instead gradually fill up ods fall under two paradigms; substitutuve and additive. the capacity of a large base model. For example, fine-tuning Substitutive methods modify an existing network to solve

consolidation, side-tuning uses fewer learnable parameters features with one or more readout layers [26]. [20, 30] mod-

updates all the parameters. A large body of constraint-based a new task by updating some or all of the network weights methods focus on how to regularize these updates in order to (simplest approach being fine-tuning). There are many alprevent inter-task interference [4]. Side-tuning does not re-ternatives that attempt to avoid catastrophic forgetting. [16] weights. The most commonly employed approaches for nettwe flact-tanking and fixed features. Given a pre-trained netquire such regularization since the additive structure means adds constraints on how the parameters are updated and [34, 16, 18] add a parameter regularization term per task. We compare side-tuning to alternative approaches on

These approaches tend to impose constraints which slow both the iCIFAR and Taskonomy datasets. iCIFAR con-down learning on later tasks (see Sec. 4.2.2 on rigidity, sists of ten distinct 10-class image classification problems. [3]). [4] relegates each task to approximately orthogonal

ulate the output by applying learned weight masks. [2, 29] 3. The Side-tuning Framework learn across different tasks with additional task-specific na-

- rameters. Perhaps the most comparable work to side-tuning to abruptly lose previously learned knowledge upon new task, learn a new network utilizing dense lateral contarget task are computed as $R(x) \triangleq B(x) \oplus S(x)$.

the base task and the current (target) task. This representation space is flexible and could either be a latent space (e.g. nature places no constraints on the structure of the side net via and add new parameters (additive learning). In incre- infants are hypothesized to learn separate, discontinuous, (see Section 4.6 for a comparison). Side-fauting is instead man a posteriori estimate and, like the MAP estimate, it. Network adaptability is the sole criterion only if we care must be detected and data cannot be replayed (e.g. to genwork, allowing the parameters to be strategically allocated. mental (lifelong) learning, substitutive methods like fine- and context-dependent perception systems [1, 17]. Adults focused on learning multiple tasks. In particular, side-tuning can use tiny networks when the tuning are at risk of forgetting early tasks. To this end, are able to rapidly learn new affordances, but only when

particular: fine-tuning, feature extraction, and other approaches are side-tuning with a fixed curriculum on the blending parameter \alpha.

tioned on one another [33].

Taskonomy covers multiple tasks of varied complexity from subspaces but is unable to transfer information across tasks. across computer vision (surface normal and depth estimation. edge detection, image 1000-way classification, etc.). lems by freezing the weights and adding a small number much smaller than the base. Consequently, even without features). One economical approach is to use off-the-shelf-

- incremental learning. Namely, it does not suffer from either:
- does not require these lateral connections, making it signif-· Rigidity: where networks become increasingly unable icantly simpler and applicable on a larger variety of probto adapt to new problems as they accrue constraints lems. Furthermore, the results suggest that side-nuring offrom previous problems. Discussed in Sec. 4.2.2. fers similar or better performance to these methods.
 - to new problems by first training on tasks sampled from could be nonparametric, and it might not be optimized for equivalent to the common (stage-wise) training curriculum Forgetting uses a decoder-based approach that can be ina standing distribution of tasks. Side-tuning is fundamenally compatible with this formulation and with existing in Section 4.4, but the simplest choice is just a pretrained that are unlocked partway through training. approaches (e.g. [8]). Moreover, recent work suggests network.



termines how heavily to weight the base model. In practice, current and previous tasks. This is the case for incremental the value of α correlates with task relevance (see Sec. 4.4). learning, where we want an agent that can learn a sequence

Figure 2. Mechanics of ride-tuning. (i) Side-tuning takes some core network (B) and adapts it to a new task by (ii) adapting a side network. (iii) Shows the connectivity structure when using side-tuning along with alpha-blending, (iv) Existing adaptation methods turn out to be special cases of side-tuning. In

this dilemma. Feature extraction ($\alpha=0$) locks the weights — networks. In principle, learning of new tasks can benefit weights are already optimal, yields a biased estimator. In nth task can use all n-1 previous tasks). Since we do not fact, the estimator allows no adaptation to new evidence make use of this available information, our results should be and is asymptotically inconsistent. On the other hand, finetuning ($\alpha = 1$), is an uninformative prior yielding a lowbias high variance estimator. With enough data fine-tuning can produce better estimates, but this usually takes more data than feature extraction. Side-tuning addresses both the these problems. It reduces variance by including the fixed features in the representation, and it is consistent because it

we solely about raw performance on a single target task. erate constraints for EWC).

Sensor Fusion Many problems are easier to solve when While the side network can be initialized using a vari-3.3. Percentual Regularization and Catastrophic While α provides a way to control the importance of the preserving performance. prior, another natural approach for enforcing a prior is to

multiplication can be though of as a measure of agreement side-tuning representation is a combination, $B(x) \oplus S(x)$. Typically, it is easier to specify meaningful explicit priors baseline for incremental learning, outperforming existing among sensors, and multiplication for fusion has been ex- What should ⊕ be? plored in [4,20]. Fil M [22] defines a combination opera. Side-tuning admits several options for this combination tations, which can be difficult if not impossible to interpret. on more tasks (in Section 4.2)

On these datasets, side-tuning uses side networks that are of new parameters per task (simplest approach being fixed in that combines between the restormance of the parameters per task (simplest approach to incre-to-in that combines between the restormance of the combines of the combin examples are concatenation and summation. We observe distance measure on the outputs can be pulled back through mental learning, which means that already-learned compothat alpha blending, $a \oplus b \triangleq \alpha a + (1 - \alpha)b$, works well the decoder and into the latent space. This induced distance nents are never undated and performance across the whole in practice. Alpha blending preserves the dimensions of the d_D on the latent representations is called the pullback met-set can only increase as the agent sees more tasks. This inputs and is simpler than concatenation. In fact, concateric in differential geometry, and in deep learning it is called monotonicity is the key property of the additive family of al-

> encommonses several common transfer learning approaches. successful application of this approach would be the auxil-As shown in Figure 2 and when the side network is the same iary losses in GPT [23], though we did not find it effective. independently of their order in the sequence (always using Base Model. The base model B(x) provides some core as the base, side-tuning is equivalent to feature extraction Perceptual regularization is often used to dampen catascognition or perception, and we put no restrictions on how when $\alpha=1$. When $\alpha=0$, side-tuning is instead equivtrophic forgetting. For example, Elastic Weight Consolidarigidity during training. We show this in Section 4.2.2. B(x) is computed. We never undate B(x), and in our apalent to fine-tuning. If we allow \(\alpha\) to vary during training tion uses a diagonalized second-order Taylor expansion of
> Side-tuning naturally bandles other continuous learning Meta-learning seeks to create agents that rapidly adapt proper it has zero learnable parameters. In general R(x) (which we generally do), then switching \(\alpha \) from 1 to 0 is the expectation of the pullback metric. Learning Without seems to 0 is the expectation of the pullback metric.

converges to the MLE (fine-tuning $\alpha = 0$)

those are minor updates to familiar, well-practiced systems [5]. On a more fine-grained scale, there are areas the problem at hand. When the base is relevant and reof functional specificity within the brain [14], including quires only a minor update, a very small network can sufwholly separate pathways where output is mutually condi-fice. Section 4.4 explores the effect of network size, how that changes with the choice of base and target tasks.

asing a suite of sensors providing complementary sources ety of methods, we initialize the side network with a copy of information [7]. We combine both a base and a resid- of the base network. When the forms of the base and side ual side network via a simple additive mechanism. This has networks differ, we initialize the side network with weights been successfully used in computer vision (ResNets [10]) distilled from the base network using knowledge distillaand in robotics, where residual RL [12, 35] learns a single tion [11]. We test alternatives in Section 4.4 control) and with a learned residual network. Alternatively, Combining Base and Side Representations. The final

on outputs (e.g. 1.2 for pixels) than on the latent represen-

Side-tuning learns a side model S(X) and combines this
nation followed by a channel-collapsing operation (e.g. 1x1 the perceptual loss [13]. This may be a useful method for
gorithms. It is worth repeating that there is No Catastroohic Catastrophic forgetting: tendency of a network is Progressive Neural Networks (PNN) [31] which, for each with a base model B(x) so that the representations for the convolution) is a strict generalization of alpha belending.

 Roowledge transfer when (i) the previous task is release to experiment in the previous task is released to experiment in the previous While simple, alpha blending is expressive enough that it the new task and (ii) there is limited training data. A recent for one of the tasks is shown in Figure 4. Furthermore.

that these approaches work primarily by feature adaptation

Side Model. Unlike the base model, the side network 0 (hyperbolic decay). In this curriculum, a controls the tion 4.2.1). Side-tuning avoids catastrophic forgetting by inflictibility becomes a serious problem for constraint-based in a series of analysis experiments. rather than rapid learning [24], and feature adaptation is also S(x) is updated during training; learning a residual that weighting of the prior (B(x)) with the learned estimate design (as the base network is never updated). we apply on top of the base encoding. Iteratively learn- (S(x)), and the weight of the evidence scales with the ing residuals for a single task is known as gradient boosting amount of data. This curriculum is suggestive of a maxi-

task by combining a coarse policy (e.g. hand-coded optimal

One crucial component of the framework is that the com-

i. Train base B(x) ii. S(x) scales with problem iii. Sidetuning

When minimizing estimation error there is often a trade- getting becomes a major issue. off between the bias and variance contributions [9]. Choos- In our experiments we dedicate one new side network to

ing between feature extraction or fine-tuning exemplifies each new task and train it independently of the earlier side and corresponds to a point-mass prior that, unless the from all the side networks learned in previous tasks (i.e. the

of tasks $T_1, ..., T_m$ and, at the end, is canable of reasonable

performance across the entire set. Thus, catastrophic for-

____ Training (Tasks) → figure 4. Theoretical learning curve of side-tuning. The model learns during task-specific training and those weights are immediately frozen,

methods and task-specific performance declines after learn-

ing more than a handful of tasks. Moreover, continuous

adaptation requires an online method as task boundaries

penalize deviations from the original feature representation. We show that this simple approach provides a strong 23. Side-tuning does not forget in incremental learning. Qualitative results for incremental learning on Taskonomy with additive learning (sidemine for 3 may) and constraint based burning (FWC bettom 3 mays). Furb may contains results for one task and columns show how medications change over the course of training. Predictions from EWC suickly deepade over time, showing that EWC still catastrophically forgets. Predictions from side-tunnot decrease, and the initial quality is better in later tasks (e.g. commare the table in surface normals). We recovide additional commarisons (including for

> keeping a small working memory of cheap side networks dom weights and trained using minibatch SGD with that constantly adapt the base network to the input task. Adam [15]. These side networks are small easy to train and when one of the networks begins performing poorly (e.g. signaling a distribution shift) that network can simply be discarded. cheap networks has found recent success in (e.g. [21]).

In the first section we show that side-tuning compares favorably to existing incremental learning approaches on both iCIFAR and the more challenging Taskonomy dataset. We mulation from [34] which scales better, giving an adthen extend to multiple domains (computer vision, RL, impullback metric. We show that such regularization does not with undefined boundaries and where there might very little itation learning, NLP) in the simplified (transfer learning) Another notable curriculum is $\alpha(N) = \frac{k}{N}$ for k > fully address the problem of catastrophic forgetting (Secscenario for N = 2 tasks. Finally, we interpret side-tunin

> We provide comparisons of side-tuning against the folproach is unable to transfer across tasks.

Parameter Superposition (PSP): A parameter-masking substitutive approach from [4] which attempts to make tasks. We first selected the twelve tasks that make predictions from side-tuning vs. EWC for a few tasks dur-

Feature extraction (features): The pretrained base network is used as-is and is not updated during training.

This is an online approach, and online adaptation with small Fine-tuning: An umbrella term that encompasses a 4.2. Incremental Learning: No Catastrophic Forety of techniques, we consider a more narrow definigetting in Additive Learning tion where pretrained weights are used as initialization

> and then training proceeds as in scratch. On both the Taskonomy dataset [36] and incremental CI- 4.2.1 Catastrophic Forgetting FAR (iCIFAR, [28]), side-tuning outperforms existing inbased substitutive approach from [16]. We use the forcremental learning approaches while using fewer param- As expected, there is no catastrophic forgetting in sideeters1. Moreover, the performance gap is larger on more tuning. Figure 5 shows that the error for side-tuning does challenging datasets. Taskanomy includes labels for multiple computer vision creases sharply for the other methods on both Taskonomy side-tune network and maintain parameter parity.

tasks including 2D (e.g. edge detection), 3D (e.g. surface and iCIFAR. normal estimation), and semantic (e.g. object classification) The difference is meaningful, and Figure 3 shows samp tasks independent from one another by mapping the tions from a single RGB image, and then created an increing and after training. As is evident from the bottom rows, weights to approximately orthogonal spaces. This apmental learning setup by selecting a random order in which EWC exhibits catastrophic forgetting on all tasks (worse to learn these tasks (starting with curvature). As images are image quality as we move right). In contrast, side-tuning

ental details (e.e. learning rate and architecture) provided proach from [31] which utilizes many lateral connec-in the supplementary.

0 1 2 3 4 5 6 7 8 9 10 11

iCIFAR. Sidetune (A) merses base and side information with a multilayer percentron (adapter).

- PNN

- FWC

Indep.

- Sidetune

- Finetune

- Features

(top) shows no forgetting and the final predictions are significantly closer to the ground-truth (boxed red).

0 1 2 3 4 5 6 7 8 9

± 60 − PNN

- Sidetune (A)

20 Sidetune

not increase after training (blue shaded region), while it in-

20 --- EWC

are 5. Incremental Learning on Taskonomy and iCIFAR. The above curves show loss and error on incremental learning experiments for three

side-tunine

ks on Taskonomy (left) and iCIFAR dataset (right). The fact that side-tuning losses are flat after training (as we go right) shows that it does not forget

recipitable learned tasks. Similarly the performance remains consistent even on later tasks (so we no down) showing that side-tasking does not become

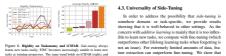
Indep.

- Finetune

0 — Features

Side-tuning learns later tasks as easily as the first while constraint-based methods such as EWC stagnate. The pre- Table 2. Average rank on Taskonomy and iCIFAR. While being redictions for later tasks such as surface normals (in Figure 3) markably simpler, side-tuning generally achieved a better average rank are significantly better using side-tuning-even immediately after training and before any forgetting can occur.

Tasknomy dataleratives.



lasknomy dataset, where side-tuning significantly outperformed all test

tuning that undated both the base and side networks: the

fine-tuning, and significantly outperformed learning from

scratch, as shown in Figure 3a. The additional structure of

the frozen base did not constrain performance with large

Habitat environment. The agent is provided with both RGB

input image and also an occupancy map of previous loca-

ture, etc.) to the the appendix. See provided code for full details.

gid. Alternative methods clearly forget (e.g. PSP) and/or become rigid (e.g. EWC). In Taskonomy, PNN and Independent are hidden under Sidetune. In the average rapidity is zero for side-noising (and almost zero for PSP). side-tuning generally performs as well as features or fine-Figure 6 quantifies this slowdown. We measure rigidity Transfer learning in Taskonomy. We trained networks as the log-ratio of the actual loss or the ith task over the to perform one of three target tasks (object classification, tions between the base and side networks. This re- 256x256 we use a ResNet-50 for the base network and a loss when that task is instead trained first in the sequence. surface normal estimation, and curvature estimation) on the quires the base model to be a neural network and the 5-layer convolutional network for the side-tuning side net-As expected, side-tuning experiences effectively zero slow- Taskonomy dataset [36] and varied the size of the training

architecture of the base and side networks to be the work. The number of learnable network parameters used down on both datasets. For EWC, the increasing constraints set $N \in \{100, 4 \times 10^6\}$. In each scenario, the base netacross all tasks is 24 6M for FWC and PSP and 11 0M for make learning new tasks increasingly difficult—and the log- work was trained (from scratch) to predict one of the non-.ndependent: Each task uses a network trained indepen-iCIFAR. First, we pretrain the base network (ResNet ratio increases with the number of tasks (Taskonomy, left). target tasks. The side network was a copy of the original dently for the target task. This method uses far more 440 on CIFAR-10. Then the 10 subsequent tasks are formed It is too rigid (log-ratio > 0) even in iCIFAR, where the later base network. We experimented with a version of finelearnable parameters than all the alternatives (e.g. sav-by partitioning CIFAR-100 classes into 10 disjoint sets of tasks are similar to earlier ones. ing a separate ResNet-50 for each task) and achieves 10-classes each. We train on each subtask for 20k steps very strong performance. Due to the scaling, it is before moving to the next one. Our state-of-the-art substi-4.2.3 Final Performance

generally not considered a viable incremental learning tutive baselines (EWC and PSP) update the base network

Overall, side-tuning significantly outperforms the other for each task (683K parameters), while side-tuning updates methods while using fewer than half the number of trainable a four layer convolutional network per task (259K parameparameters of the other methods. When the other methods amounts of data (4M images), and side-tuning performed use smaller networks, their performance decreases further. as well as (and sometimes slightly better than) fine-tuning. On both iCIFAR and Taskonomy, side-tuning achieves the

Ouestion-Answering in SQuAD v2. We also evaluated best is 1.88 (PNN) see Fig. 2)

mains). In fact, the much larger networks used in EWC

and PSP should achieve better performance on any single

task. For example, EWC produces sharper images early on

in training, before it has had a chance to accumulate too

many constraints (e.g. reshading in Fig. 3). But this factor

was outweighed by side-tuning's immunity from the effects

of catastrophic forgetting and creeping rigidity.

best average rank (1.13 of 4 on Taskonomy, while the next side-tuning on a question-answering task (SQuAD v2 [25]) using a non-convolutional architecture. We use a pretrained This is a direct result of the fact (shown above) that side. RERT [6] model for our base and a second for the side tuning does not suffer from catastrophic forgetting or rigid-networkdefer. Unlike in the previous experiments, BERT ity. It is not due to the fact that the sidetuning structure uses attention and no convolutions. Still, side-tuning adapts is specially designed for these types of image tasks; it is to the new task just as well as fine-tuning, outperforming not (we show in Sec. 43 that it performs well on other dofeatures and scratch (Figure 3h)

tunine-whichever is better

Imitation Learning for Navigation in Habitat, Wo In order to test whether side-tuning could profitably syntrained an agent to navigate to a target coordinate in the

> using 49, 490, 4900, or 49k expert trajectories and pretrained denoising features. Side-tuning was always the best-

iCIFAR (Figure 8 Left), if catastrophic forgetting is not a concern then the parameters would've been better used in a deeper network rather than many shallow networks. Table 4. Average rank of various merge methods. There are multiple orming as well or better than the best competing method in each domain: (a) In Taskonomy, performing either Normal Estimation or yield a minor boost in performance. We found that

biject Classification using a base trained for Curvanues and either 100 or 4M images for transfer. Results using Obj. Cls. base are similar initializing from the base network slightly outperforms and provided in the amendix. (b) In SOuAD v2 question-answering, using BERT instead of a convolutional architecture. (c) In Habitat. learning to navigate by imitating expert navigation policies, using inputs based on either Curvature or Denoising. Finetuning does not

detection and autoencoding) and find the difference in initialization is now significant (p = 0.01). one from the side model. Are both streams necessary? Fig-More than just stable updates. In RL, fine-tuning often fails to improve performance. One common rational- improves when using both models. ization is that the early updates in RL are 'high variance'. The usual solution is to first train using fixed features and 5. Discussion then unfreeze the weights at some point in training (via a

hyperparameter to be set). We found that this stage-wise We have introduced the side-tuning framework, a simple approach performs as well (but no better than) keeping the yet effective approach for additive learning. Since it does 7. Side-tuning outperformed alternatives on intermediat features fixed-and side-tuning performed as well as both not suffer from catastrophic forgetting or rigidity, it is natu-

interaction instead of expert trajectories, we observe identical trends. We trained agents directly in Habitat (74 builds (e.g. 4.9k trajectories) outperformed the other technique ings). Fig. 3d shows performance in 14 held-out buildings (side-tune 9.3 vs. fine-tune: 7.5, features: 6.7, scratch: 6.6 after 10M frames of training. Side-tuning performs compa- Figure 7). Network size. Does network size matter? We find (i) If

the target problem benefits from a large network (e.g. clas-

sification tasks), then performance is sensitive to side net-

rably to the max of competing approaches. results were similar to standard fine-tuning 3. In all scenar-4.4. Learning Mechanics in Side-Tuning ios, side-tuning successfully matched the adaptiveness of

Task relevance predicts alpha α . In our experiments, work size but not size of the base. (ii) The base network we treat α as a learnable parameter and find that the relative values of α are predictive of emprical performance. In will still offer advantages over alternatives. In the suppleimitation learning (Fig. 4.3), curvature ($\alpha = 0.557$) out mentary material we provide supporting experiments from performed denoising ($\alpha=0.252$). In Taskonomy, the α Taskonomy using both high- and low-data settings (curvavalues from training on just 100 images predicted the ac- ture → {obj. class, normals}, obj. class → normals), and in tual transfer performance to normals in [36], (e.g. curvature Habitat (RL using {curvature, denoise} → navigation). ($\alpha = 0.56$) outperformed object classification ($\alpha = 0.50$)).

perform as well in this domain. (d) Using RL (PPO) and direct interaction instead of supervised learning.

tions. The man does not contain any information about the

environment-just previous locations. In this section we

use Behavior Cloning to train an agent to imitate experts

following the shortest path on 49k trajectories in 72 build-

ines. The agents are evaluated in 14 held-out validation

buildings. Depending on the what the base network was

trained on, the source task might be useful (Curvature) or

harmful (Denoising) for imitating the expert and this deter-

mines whether features or learning from scratch performs

best. Figure 3c shows that regardless of the which approach

Reinforcement Learning for Navigation in Habitat.

For small datasets, usually $\alpha \approx 0.5$ and the relative order

rather than the actual value is important.

sing a different learning algorithm (PPO) and using direct

worked best, side-tuning consistently matched or beat it.

Benefits for intermediate amounts of data. We showed in the previous section that side-tuning performs like the best of {features, fine-tuning, scratch} in domains with abundant or scant data.

> thesize the features with intermediate amounts of data, Figure 8. Boosting and Gradient Variance. (Left) Deeper network we evaluated each approach's ability to learn to navigate perform many shallow learners. (Right) Features and Sule-tuning do more

Features 0.204/0.117 24.4/45.4 49.4 49.5 11.2 8.2 11.9

performing approach and, on intermediate amounts of data Not Boosting. Since the side network learns a resid-

could glean by extending side-tuning to do boosting? Although network boosting this does improve performance on 3. Side-tuning comparisons in other domains. Sidetuning matched the adaptability of fine-tuning on large datasets, while per-

Initialization. A good side network initialization can tones of feature, wise transformation to choose from. We decided to nick a low-energy initialization, which slightly outperforms Xavier initialization. However, we found that these differ-

while being simpler than stage-wise (Fig. 8 Right). We rally suited to incremental learning. The theoretical advan-

tested the 'high-variance update' theory by fine-tuning with tages are reflected in empirical results, and side-tuning out-

both gradient clipping and an optimizer designed to prevent performs existing approaches in challenging contexts and

such high-variance updates by adaptively warming up the with state-of-the-art neural networks. We further demon-

two streams of information - one from the base model and tasks are presented in a sequence and that task identities are

ences were not statistically significant across tasks (Ho : pretrained = xavier; p = 0.07, Wilcoxon signed-rank test). We suspect that initialization might be more important on Table 5. Average rank on Taskonomy when ablating base and side models. Performance improved when using both the base and side modharder problems. We test this by repeating the analysis without the simple texture-based tasks (2D keypoint + edge

Merge methods. Section 2 and 3.1 described differ-

evaluates a few of these approaches on the Taskonomy 6. Limitations

ent ways to merge the base and side networks. Figure 4

dataset. Element-wise product and addition via alpha-

blending are two of the simplest approaches and have lit-

tle overhead in terms of compute and parameter count. [31]

tron (MLP). FiLM [22] adds additional compute by defin-

ing $a \oplus b = \gamma_{\theta}(a) \odot b + \beta_{\theta}(a)$ where γ_{θ} and β_{θ} is an MLP

with two heads. Figure 4 shows that all three approaches

are roughly comparable, though the MLP-based methods

achieve marginally better average rank on the Taskonomy

dataset. Nonetheless, we chose to use the simplest feature-

wise transformation (alpha-blending) since it adds no pa-

re the side network is trained so that it does not impact the output.

rameters and achieves similar performance.

Full details in the supplementary.

uses $a \oplus b = F_0(a) + b$ where F_0 is an a Multi-Layer Percep-

ual on top of the base network, we ask: what benefits we

lifelong-learning approaches, while being significantly sim-

could be analyzed and subsequently relaxed. In particular:

cremental learning experiments used the same side network

architecture for all subtasks. A method for automatically

adapting the networks to the subtask at hand could make

more efficient use of the computation and supervision. Bet-

ter forward transfer: Our experiments used only a single

base and single side network. Leveraging the already previ-

known. Using several active side networks in tandem would

ition (Alpha Blending)

[3] Arslan Chaudhry Manc' Aurelio Ranzato, Marcus Robrbach GEM. CoRR, abs/1812.00420, 2018. 2

velopmental science, 14(2):306-318, 2011.

any other Toyota entity.

and Bruno A. Olshausen. Superposition of many models into one CoRR abs/1902-05522 2019 2 3 5

CoRP abr/1901 06510 2019 2 3

tuning are not solely due to gradient stabilization early in tuning to perform on-par-with or better-than many current formers for language understanding. CoRR, abs/181004805.

CoRR, abs/1703.03400, 2017. 2

tion, 4(1):1-58, 1992, 4

recognition, pages 770-778, 2016. 3 11) Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distill-

ously trained side networks could yield better performance on later tasks. Learning when to deploy side networks: arXiv-1503 02531 2015 3

Ablating Base and Side Elements. Side-tuning uses Like most incremental learning setups, we assume that the [12] Tobias Johannink, Shikhar Bahl, Ashvin Nair, Jianlan Luo,

tracked task relevance, but a more rigorous treatment of the tual losses for real-time style transfer and super-resolution interaction between the base network, side network, \(\alpha \) and \(\cong \) CoRR, abs/1603.08155, 2016. 4 final performance could yield insight into how tasks relate [14] Nancy Kanwisher. Functional specificity in the human

Acknowledgements: This material is based upon work 107(25):11163-11170, 2010, 3 supported by ONR MURI (N00014-14-1-0671), Google [15] Diederik P Kingma and Jimmy Ba. Adam: A method for Cloud. NSF (HS-1763268), NVIDIA NGC beta, and TRI. stochastic optimization. International Conference on Learning Representations, 2015. 5

the opinions and conclusions of its authors and not TRI or Joel Veness Guillaume Designins Andrei A Rusu Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Grabsk

opmental continuity? crawling, cruising, and walking. De-[17] Kari S Kretch and Karen E Adolph. Cliff or step? posturespecific learning at the edge of a drop-off. Child develop

The missing link between faces, text, planktons, and cat breeds, arXiv preprint arXiv:1701.07275, 2017. 2 CoRR. abs/1606.09282, 2016, 2

and Mohamed Elhoseiny. Efficient lifelong learning with A [4] Brian Cheung, Alex Terekhov, Yubei Chen, Pulkit Agrawal.

[1] Karen E Adolph, Sarah E Berger, and Andrew J Leo. Devel-

5) Whitney G Cole, Gladys LY Chan, Beatrix Vereilken, and

Karen F Adolph Perceiving affordances for different motor skills. Experimental brain research, 225(3):309-319, 2013. manan, and Kayvon Fatahalian. Online model distillation

learning rate (RAdam, [19]). This provided no benefits strated that the approach is effective in multiple domains [6] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina over vanilla fine-tuning, suggesting that the benefits of side-

moulin, and Aaron Courville. Film: Visual reasoning with [7] Vincent Dumoulin, Ethan Perez, Nathan Schucher, Florian Strub. Harm de Vries. Aaron Courville, and Yoshua

Rengio Feature-wise transformations Distill 2018 https://distill.pub/2018/feature-wise-transformations. 3 The naïve approach to incremental learning used in this [8] Chelsea Finn, Pieter Abbeel, and Sergey Levine. Modelpaper made a number of design decisions. These decisions agnostic meta-learning for fast adaptation of deep networks

Flexible parameterizations for side networks: Our in- [9] Stuart Geman, Elie Bienenstock, and René Doursat. Neural networks and the bias/variance dilemma. Neural computa-

101 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In Proceedings of the IEEE conference on computer vision and pattern [25] Pranav Rajpurkar, Robin Jia, and Percy Liang. Know what

brain: a window into the functional architecture of the mind. Proceedings of the National Academy of Sciences

Toyota Research Institute ("TRI") provided funds to assist the authors with their research but this article solely reflects [16] James Kirkpatrick, Razyan Pascanu, Neil C. Rabinowitz.

ing side-tuning to measure task relevance: We noted that \(\text{ } \) [13] Justin Johnson, Alexandre Alahi, and Fei-Fei Li. Percep-

Barwinska, Demis Hassabis, Claudia Clopath, Dharshan K maran, and Raia Hadsell. Overcoming catastrophic forge ting in neural networks. CoRR, abs/1612.00796, 2016.

[2] Hakan Bilen and Andrea Vedaldi. Universal representations: ment 84(1):226-240 2013 2 [18] Zhizhong Li and Derek Hojem, Learning without forgetting

[19] Liyuan Liu Haoming Jiang Penocheng He Weizhu Che

Xiaodong Liu, Jianfeng Gao, and Jiawei Han. On the variance of the adaptive learning rate and beyond, arXiv preprint arXiv:1908.03265, 2019, 9

[20] Arun Mallya and Svetlana Lazebnik. Piggyback: Adding multiple tasks to a single, fixed network by learning to mask

[21] Ravi Teia Mullanudi. Steven Chen. Kevi Zhang. Deva R.

for efficient video inference. CoRR, abs/1812.02699, 201 [22] Ethan Perez, Florian Strub, Harm De Vries, Vincent I

general conditioning layer. In Thirty-Second AAAI Confe ence on Artificial Intelligence, 2018, 3, 9 [23] Alec Radford, Karthik Narasimhan, Tim S

mans, and Bya Sutskever. Improving langu understanding by generative pre-training. https://s3-us-west-2. amazonaws. com/openai assets/researchcovers/languageunsupervised/language understanding naner pdf 2018 4

[24] Aniruddh Raghu, Maithra Raghu, Samy Bengio, and Oriol Vinvals, Rapid Learning or Feature Reuse? Towards Unde standing the Effectiveness of MAML. arXiv e-prints, page arXiv:1909.09157. Sep 2019. 2

you don't know: Unanswerable questions for squad. CoRR abs/1806.03822, 2018, 2, 7

ing the knowledge in a neural network. arXiv preprint [26] Ali Sharif Razavian, Hossein Azizpour, Josephine Sullivan, and Stefan Carlsson. CNN features off-the-shelf: an astound

Avinash Kumar, Marthias Loskyll, Juan Aparicio Ojea, Eu- [27] Sylvestre-Alvise Rebuffi. Alexander Kolesnikov. an gen Solowiow, and Sergey Levine. Residual reinforcement Christoph H. Lampert, icarl: Incremental classifier at

ing baseline for recognition. CoRR, abs/1403.6382, 2014.

provide a natural way to detecting distribution shift. Us- learning for robot control. CoRR, abs/1812.03201, 2018. 3 representation learning. CoRR, abs/1611.07725, 2016. 2