# Model Stitching by Functional Latent Alignment

**Ioannis Athanasiadis    Anmar Karmush    Michael Felsberg**
Linköping University
{firstname.lastname}@liu.se

## Abstract

Evaluating functional similarity involves quantifying the degree to which independently trained neural networks learn functionally similar representations. Reliably inferring the functional similarity of these networks remains an open problem with far-reaching implications for AI. Model stitching has emerged as a promising paradigm, where an optimal affine transformation aligns two models to solve a task, with the stitched model serving as a proxy for functional similarity. In this work, we draw inspiration from the knowledge distillation literature and propose Functional Latent Alignment (FuLA) as a novel optimality condition for model stitching. We revisit previously explored functional similarity testbeds and introduce a new one, based on which FuLA emerges as an overall more reliable method of functional similarity. Specifically, our experiments in (a) adversarial training, (b) shortcut training and, (c) cross-layer stitching, reveal that FuLA is less prone to artifacts tied to training on task cues while achieving non-trivial alignments that are missed by stitch-level matching.

## 1    Introduction

Deep neural networks lie at the forefront of modern AI, driving significant advances across various domains ranging from computer vision [24] and natural language processing [31] to healthcare [14]. Their success is attributed to their ability to learn meaningful representations of the data [4, 7] which capture their abstract relationships [13, 45]. Understanding the emergence of these representations constitutes an important problem [23] from both scientific and applied perspectives [12]. One approach to gaining insight into learned representations is to compare the similarity of internal representations of different models. In the literature, such similarities generally fall into two categories: representational or functional similarity. In this work, we focus on the latter and address the problem of evaluating functional similarity. We adopt a functional consistency view: functionally similar networks should exhibit similarity in their internal transformations leading to the output.

A plethora of representational similarity metrics have been proposed such as projection weighted Canonical Correlation Analysis (PWCCA) [33], Singular Vector Canonical Correlation Analysis (SVCCA) [33] and the wide-spread Centered Kernel Alignment (CKA) [23]. These structure-based metrics measure correlation between geometric structures in activation spaces of independently trained models given the same data. Although they have been frequently used in the literature to extract relevant insights such as in [30, 32, 36], they have also been criticized for overestimating similarity [21] due to spurious feature correlations, as well as for being agnostic to the functional behavior [10] and invariance properties [35] of the model.

Model stitching [3, 9] is a prominent paradigm for evaluating the functional similarity of neural networks. In contrast to the structure-based metrics, the similarity is grounded on whether the representations of one model can be linearly aligned to those of another one in such a way that a functional property of interest (e.g., accuracy) is maintained. In practice, this is realized by freezing the two halves of these neural networks and connecting them by a trainable affine transformation which is responsible for aligning the models.
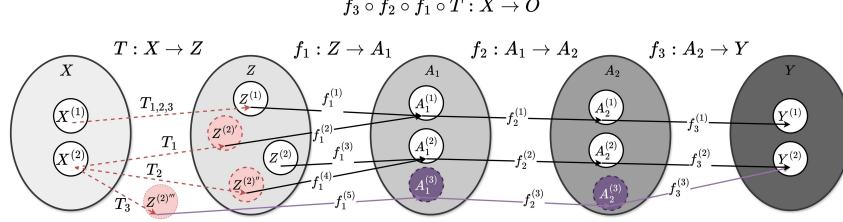
Figure 1: Conceptual visualization for different degrees of functional alignment. The $f_3 \circ f_2 \circ f_1$ is composition of functions mapping the input domain $\boldsymbol{Z}$ into the output domain $\boldsymbol{Y}$ through two intermediate domains $A_1$ and $A_2$. Let $T_1$, $T_2$ and $T_3$ be three different transformations, mapping domain $\boldsymbol{X}$ to $\boldsymbol{Z}$, that only differ in how they map the subdomain $\boldsymbol{X}^{(2)}$ to $\boldsymbol{Z}^{(2)'}$, $\boldsymbol{Z}^{(2)''}$ and $\boldsymbol{Z}^{(2)'''}$ respectively – each equidistant from $\boldsymbol{Z}_2$. Under the DM, all three transformations are treated as equivalent. In contrast, the task-oriented objectives regards $T_1$ as the least optimal, while considering $T_2$ and $T_3$ equally performant, since they both preserve the input–output relationship of the composition.

Although less popular than the structure-based metrics, model stitching has been used to draw insights [3, 9], such as demonstrating that non-robust and robust models are functionally similar with respect to clean accuracy [1]. However, recent works questioned the reliability of model stitching under the commonly used setting - task loss matching (TLM) [2, 19, 40]. In response to that, direct matching (DM), linearly connecting neural networks by minimizing the representation distance at the stitching level - was advocated as a more reliable alternative for evaluating similarity.

In this work, we take a fresh look at model stitching and argue that a reliable functional similarity metric should reflect the degree of functional comparability between two models trained under different settings, where trivial solutions may not exist (e.g., different datasets and/or different training recipes). In this context, the affine transformation used for representation alignment should both (i) effectively propagate functionally relevant information and (ii) avoid compensating for functional discrepancies between the models.

Based on these two criteria, we hypothesize that the DM is at risk of violating (i) as it does not account for the functional nuances of the layers following the stitching (e.g., $T_1$ transformation in Fig 1). Conversely, the TLM might fail in regards to (ii) by exploiting irregularities in the decision process of the model to achieve optimal task performance (e.g., $T_3$ transformation in Fig. 1).

Towards developing model stitching that better fits the needs of functional similarity evaluation, we draw inspiration from the work conducted by Liu et al. [27] in which the notion of function-consistency was employed for Knowledge Distillation (KD). We propose model stitching by Functional Latent Alignment (FuLA) as a novel setting for functional similarity evaluation in which the affine transformation minimizes not only the feature representations at the stitching level (i.e., DM) but also at the layers following it (see Fig. 2). Given that the alignment transformation is the only trainable layer in model stitching, functional consistency can be achieved by simply aligning the latent feature representations. Additionally, under FuLA only the representations between the stitching and the penultimate layer are aligned, that is, ignoring the output level cues. The motivation behind this choice was to avoid exploiting irregularities in the decision boundary through leakage of task-specific information during the alignment. Concretely our contributions are:

- We propose FuLA as a novel setting for model stitching under which function-consistent alignment is realized to infer the functional similarity of two neural networks.
- We show that task-based stitching can fabricate alignment by overfitting on the task.
- Through targeted experiments, we demonstrate that FuLA is less prone to these pitfalls and captures non-trivial alignments missed by direct matching, making it a more reliable metric for functional similarity.

## 2  Model Stitching

In this section, we introduce the notation and provide an overview of the components that form the basis of our experimental setup, finally we introduce FuLA as a model stitching setting.
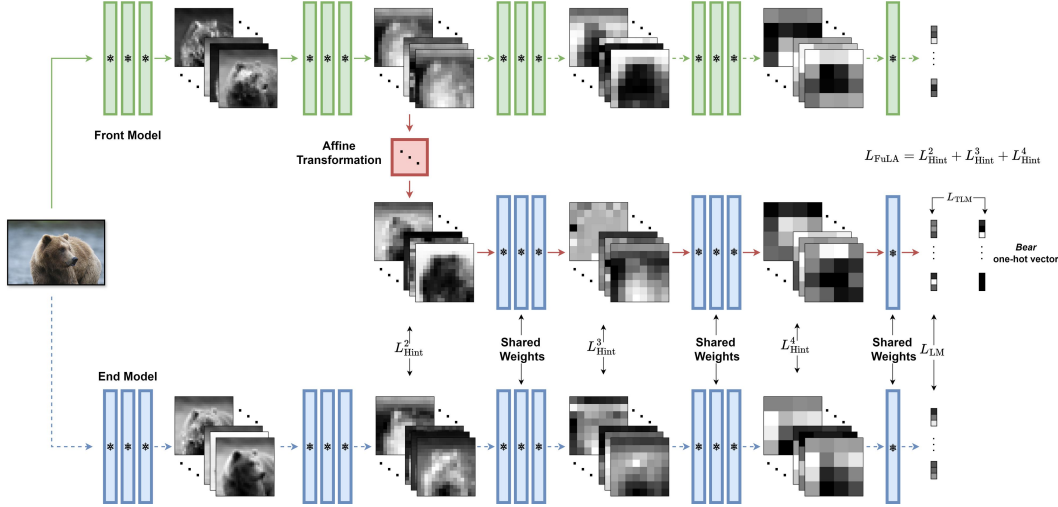
Figure 2: The proposed model stitching by FuLA in relation to other model stitching settings. In the example of the figure, stitching is performed at the second layers between identical architectures.

## 2.1 Notation

We follow the notation used by Balogh and Jelasity [2]. Let $f_w : \mathcal{X} \to \mathcal{Y}$ be a feedforward neural network with $m$ layers, parameterized by $w$ and realized as the composition $f = f_m \circ \cdots \circ f_1$, where each $f_i : \mathcal{A}_{f,i-1} \to \mathcal{A}_{f,i}$ maps the activation space of the $(i-1)^{\text{th}}$ to that of the $i^{\text{th}}$, with $\mathcal{A}_{f,0} = \mathcal{X}$. For notational simplicity, we omit the explicit dependence on $w$ when clear from context. We define the partial compositions $f_{>i} = f_m \circ \cdots f_{i+1}$ and $f_{\leq i} = f_i \circ \cdots f_1$ such that $f = f_{>i} \circ f_{\leq i}$.

Given two neural networks $f$ and $g$ with $m$ and $k$ layers, respectively, model stitching aims to quantify the extent to which $g_{>j}$ can replicate the functionality of $g$ when given as input the $\mathcal{A}_{f,i}$ with $i \in [m]$ and $j \in [k]$. Formally, the stitched model is constructed as $h = g_{>j} \circ T_\theta \circ f_{\leq i}$. Here the $T_\theta : \mathcal{A}_{f,i} \to \mathcal{A}_{g,j}$ is referred to as stitching layer, and we refer to $f$ and $g$ as front and end models, respectively. In practice, the optimal $T_\theta$, under a relevant optimality criterion, is selected from a suitable family of transformations $S$, while $g$ and $f$ are kept fixed. Given the nature of the problem, the transformation $T_\theta \in S$ needs to be sufficiently flexible to allow non-trivial mappings while ensuring that the capacity of the $h$ does not exceed that of the combined partial networks $g_{>j}$ and $f_{\leq i}$. The family $S$ of affine transformations is considered in the literature to satisfy these requirements [2].

We assume access to two datasets, $D_{\text{train}} = \{(x_i, y_i)\}_{i=1}^n$ and $D_{\text{test}} = \{(x_i, y_i)\}_{i=1}^c$, drawn from the underlying distributions $p(x,y)$ and $q(x,y)$, respectively. The distributions $p$ and $q$ may be identical or distinct. The $T_\theta$ is optimized using $D_{\text{train}}$, while the functional similarity is evaluated on the $D_{\text{test}}$.

**Soft label matching (SLM):** Under SLM [9], the output of the end model is used as soft label and the $T_\theta$ is obtained as:

$$\arg \min_\theta \mathbb{E}_{p(x,y)}[\mathcal{L}_{\text{SLM}}((g_{>j} \circ T_\theta \circ f_{\leq i})(x), g(x))], \text{ with } L_{\text{SLM}} : \mathcal{Y} \times \mathcal{Y} \to \mathbb{R}. \qquad (1)$$

**Task loss matching (TLM) [1]:** Under the TLM [3, 9], $T_\theta$ is found as:

$$\arg \min_\theta \mathbb{E}_{p(x,y)}[\mathcal{L}_{\text{TLM}}((g_{>j} \circ T_\theta \circ f_{\leq i})(x), y)], \text{ with } L_{\text{TLM}} : \mathcal{Y} \times \mathcal{Y} \to \mathbb{R}. \qquad (2)$$

For example, in the classification setting we considered, both $L_{\text{SLM}}$ and $L_{\text{TLM}}$ are the cross-entropy losses between the output of the stitched model and the soft labels (from the end model) and the ground-truth (hard) labels, respectively.

---

[1]Csiszárik et al. [9] use the term TLM interchangeably to refer to both TLM and SLM, and argue that these two settings lead to highly correlated results. However, in our work, we find that TLM and SLM can behave differently under certain key stitching settings, and therefore, we distinguish between them.

**Direct matching (DM):** Under the DM, the objective for the $T_\theta$ is task-agnostic and is defined as:

$$\arg \min_\theta \mathbb{E}_{p(x)} \left[ || \left( T_\theta \circ f_{\leq i} \right) (x) - g_{\leq j}(x) ||_F \right], \tag{3}$$

As discussed in [9], under DM, there exists a closed-form solution for $T_\theta$, making it well-suited for initialization prior to further fine-tuning under SLM or TLM. Moreover, model stitching, under DM, has recently been advocated as a more reliable alternative to TLM [2], in light of the seemingly counterintuitive behavior related to the directionality of functional similarity in model stitching.

## 2.2 Adversarial Training

Adversarial Training (AT) [29] is widely regarded as the most effective method for improving model robustness against worst-case perturbations, referred to as adversarial examples [17, 26]. Functional similarity between non-robust and robust models was investigated under TLM in [1] using the AT strategy proposed in [29]. In this case, the training objective for the neural network $f_w$ becomes:

$$\arg \min_w \mathbb{E}_{p(x,y)} \left[ (1 - \alpha)\mathcal{L}(f_w(x), y) + \alpha \max_{\delta \in \mathcal{B}(x,\epsilon)} \mathcal{L}(f_w(\mathbf{x} + \delta), y) \right], \tag{4}$$

where the parameter $\alpha \in [0, 1]$ controls the proportion of adversarial samples, $\mathcal{B}(x, \epsilon)$ denotes the $\ell_\infty$-ball of radius $\epsilon$ and $\mathcal{L}$ the loss used to train the benign model $f_w$. In this study, we revisit functional similarity of robust and non-robust model through the lens of functional alignment.

## 2.3 Functional Latent Alignment (FuLA)

When evaluating functional similarity under model stitching, it is widely accepted that the family of transformations used to connect the two models, must be of low capacity [3, 9, 25], to ensure that it does not compensate for potential misalignment while solving the optimization objective. In our work, we expand upon this notion and argue that it is the synergy between the transformation family and the optimization objective that must adhere to the qualities expected of a functional similarity metric. Our view is that functional similarity should reflect the extent to which two models share their underlying functional properties when transforming input into output. In this regard, the existing model stitching settings, for functional similarity, do not account for (miss)alignment in the internal processes of the models, which can impair their reliability as metrics. To this end, we propose model stitching by FuLA where the stitched model is explicitly optimized for internal functional alignment in a task-agnostic manner. An overview of the method and its relation to existing stitching settings in provided in Fig. 2. In particular, under FuLA, the alignment transformation $T_\theta$ of the stitched model $h = g_{>j} \circ T_\theta \circ f_{\leq i}$ is found as:

$$\arg \min_\theta \mathbb{E}_{p(x)} \left[ C_j \underbrace{\frac{|| \left( T_\theta \circ f_{\leq i} \right) (x) - g_{\leq j}(x) ||_F}{||g_{\leq j}(x)||_F}}_{L_{\text{Hint}}^j} + \sum_{l=j+1}^{k} C_l \underbrace{\frac{||g_{\leq l}((T_\theta \circ f_{\leq i}) (x)) - g_{\leq l}(x)||_F}{||g_{\leq l}(x)||_F}}_{L_{\text{Hint}}^l} \right], \tag{5}$$

where $C \in \{c \in \mathbb{R}^{k+1-j} : c_l \geq 0, \sum_{l=1}^{k+1-j} c_l = 1\}$ defines a weighted average that controls the relative contribution of each term. We use uniform weighting as the default option. Additionally, we divide by the target's norm to account for potential differences in scale between feature activations at different depths. We refer to these terms as Hints [39], denoted as $L_{\text{Hint}}^t$, where $t$ indicates the depth. Note that for $t = j$, the structural Hint $L_{\text{Hint}}^j$ corresponds (up to scale) to the direct matching objective defined in Eq (3), while the functional Hints $L_{\text{Hint}}^l$ for $t = l \geq j + 1$ realize function alignment. Namely, the transformation $T_\theta$ is trained to map the $A_{f,i}$ into $A_{h,i}$ such that the latter is interpreted by $g_{>j}$ similarly to its native input $A_{g,j}$ originating from $g_{\leq j}$.

# 3 Experiments

We hypothesize that commonly used model stitching settings, often employed as proxies for assessing functional similarity, may be unreliable for this purpose. Since ground truth for functional similarity is inherently unavailable, we design a series of controlled experiments to expose the limitations of existing approaches and evaluate how FuLA compares against them.

**Experimental setting:** We consider image classification networks and conduct our experiments using CIFAR-10 and a low-resolution version of ImageNet [11], from which we randomly sampled 10 classes. Throughout the paper, we use the ResNet18 [18] architecture as our trained front and end models. We perform stitching at the first convolutional layer and at each residual residual block. In the supplementary material, we provide additional implementation details and results on more architectures, where the conclusions of our paper remain unchanged.

**Training of stitching layer:** When training the stitching layer $T_\theta$ we follow [9] and model it as a $1 \times 1$ convolutional layer with bias. The front and end models are frozen where only their running batch statistics (RS) are updated. The stitching layer is initialized by linearizing the DM objective in Eq. (3), then solved using the Moore-Penrose pseudoinverse with 100 training samples [2, 9].

**The stitching plot:** When stitching corresponding layers (i.e., $i = j$), we use the stitching plot [1], where the y-axis represents a relevant performance metric and the x-axis denotes the depth of the front layer composition. The end points correspond to the baseline end and front models.

**Implementing DM as Hint:** To better understand how different optimization settings in model stitching impact functional similarity evaluation, we focus on the relevant aspects of this comparison. Based on our analysis, we found that DM, as proposed by Balogh and Jelasity [2], may be unfairly underperforming compared to SLM and TLM due to two key factors: (i) using fewer samples to learn $T_\theta$ and (ii) not accounting for the running batch statistics of the training set, as shown in Fig. 3. We address this by learning the $T_\theta$ under the DM objective in Eq. (3) using the full training set (i.e., implemented as structural Hint [39]) [2], similar to the SLM, TLM and FuLA settings. We clarify that we refer to direct matching as DM only when it is used for initializing the stitched model, and as a Hint otherwise.
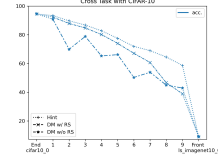


Figure 3: DM in practice.

**Experimental Outline and Takeaways:** First, we show in 3.1 that task-based model stitching can fabricate functional alignment depending on the training objective, which undermines its reliability as a functional similarity metric when transitioning between models that differ in or lack necessary functional properties. We reach a similar conclusion in 3.2, where task-based stitching is shown to (i) base alignment on visual cues unseen by both models or (ii) overlook alignment when not required by the task, both of which hinder the reliability. In light of these findings, in 3.3 we revisit cross-layer stitching and question whether high functional similarity between distant layers is an artifact of task-based alignment or an inherent property of the models' internal functional behavior, and conclude the latter. Overall, we find that FuLA, a task-agnostic method, retrieves non-trivial transformations missed by stitch-level alignment and is significantly less susceptible to the pitfalls of task-based approaches.

## 3.1 Model Stitching under Adversarial Training

We start by revisiting functional similarity between non-robust and robust models [1] as its multi-objective nature will allow us to gain insights on the behavior of different model stitching settings. In this case, the functional similarity can be evaluated based on two key and competing [41] functional properties, namely accuracy and robust accuracy [1]. We train the robust models using AT according to (4) with $\alpha = 0.5$. At all times, we train the stitching layer on the CIFAR-10 and experiment with same- and cross-task settings.

In Figs. 4 and 5 we compared model stitching under different setting for various configurations of adversarial examples proportions, front and end models. When looking at model stitching without AT (i.e. $\alpha = 0$) in the first row of Fig 4, we observe that all settings interpolate the performance of their models which aligns with what was reported originally by Balogh and Jelasity [1].

---

[2]This choice comes with an additional computational cost compared to the single-pass pseudoinverse solution, but it is necessary for disambiguating the effects of the different stitching settings.
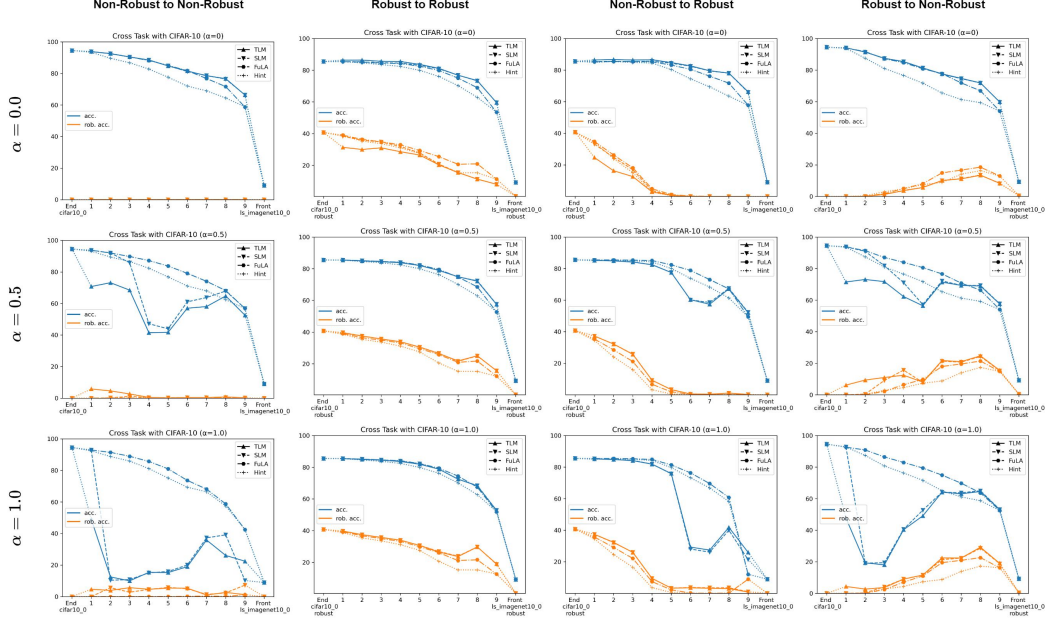
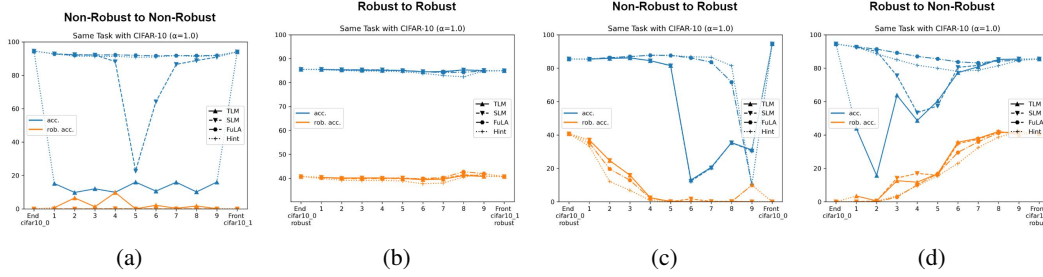Figure 4: Cross-task stitching under AT – ResNet18.



Figure 5: Same-task stitching under AT ($\alpha = 1$) – ResNet18

Model stitching with AT poses particular interest, as distinct behaviors emerge across different stitching configurations (e.g., first column in Fig. 4). When stitching non-robust models, the SLM and TLM setting display an irregular behavior in both clean and robust accuracy, which becomes more pronounced when only adversarial examples are used (i.e. $\alpha = 1$). In particular, these settings hallucinate non-robust accuracy at the cost of clean accuracy. This is a highly problematic behavior, as functional similarity metrics should capture genuine functional comparability rather than fabricate it to improve task performance. Similar, irregular behavior can be observed when stitching robust to non-robust and the reverse, where clean accuracy drops in an attempt to achieve robust accuracy earlier or for longer depending on the exact configuration.

On the other hand, the Hint and FuLA behave more stably across all cases, as they smoothly interpolate the functional properties between the two extremes. Overall, FuLA tends to produce better-performing stitched models without exhibiting signs of fabricated alignment. A reoccurring exception was the non-robust to robust configuration where in most cases Hint reports higher accuracy in the deeper layers when training using $\alpha = 1$. Interestingly, this is also the configuration where both Hint and FuLA consistently report non-existent functional similarity (i.e., random chance performance) at the penultimate layer (i.e., the $9^{th}$ layer corresponds to the layer before the fully connected classifier).

We further investigated this behavior by analyzing the CKA [23] and the functional similarity at the penultimate layers of the end and the stitched models. For the CKA, we used a DM-initialized stitched model and evaluated the similarity between its end model, using adversarial examples crafted on the former. For functional similarity, we used model stitching under FuLA by setting the contributions of

all (functional) Hints weights to 0 apart from the last which was set 1 (see Eq. (5)). Fig. 6 suggests a distinct functional behavior as we morph from the robust to non-robust models.

In particular, the penultimate layers exhibit a decaying trend in functional similarity for the non-robust to robust stitching configuration, starting at $6^{th}$ layer, whereas the others remain relatively steady throughout all layers. In light of this, FuLA may more accurately detect the onset of functionally dissimilar representations earlier than Hint (see Fig. 5c), as per its use of functional Hints.

Figure 6: Adversarial example penultimate layer similarity.

While CKA also identifies the penultimate layers between the robust model and the non-robust-to-robust model stitched at $9^{th}$ as the least similar, it fails to capture the functional nuances meaningfully. More specifically, the relative differences in CKA do not appear to be predictive of the functional comparability of the front and end models at their penultimate layers. Similar limitations of CKA have also been reported in [10, 35].

**Takeaway:** Task-based model stitching can be unreliable as a functional similarity metric fabricating alignment to solve the task. In contrast, alignment-based methods behave more intuitively, with FuLA generally achieving higher metrics without the drawbacks of optimizing for task performance.

### 3.2 Stitching under Shortcuts

Next, we explore a setting where shortcuts[16], previously unseen by both the front and end models, are introduced during model stitching. We adapt the AT stitching setup for shortcut learning using models trained on clean data and stitching with three types of shortcuts: pattern-based, location-based, and their combination. Each class is assigned a unique $4 \times 4$ marker—either a single-colored square (pattern-based) or uniformly random noise at a specific location (location-based) (see Fig. 7).
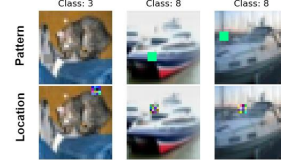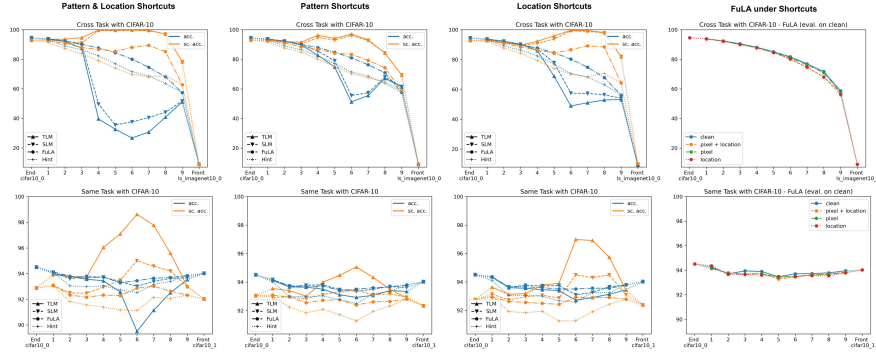
Figure 7: Pattern-based vs location-based shortcuts.

Figure 8: Stitching under shortcuts.

In Fig 8, we report the functional similarity across the three key shortcut conditions considered, for the different model stitching settings. We note that shortcut reliance can manifest in two primary ways: (i) improved performance on shortcut-augmented data relative to clean data, and/or (ii) degraded performance on clean data due to the model's overreliance on spurious cues. Both (i) and (ii) are problematic as they can provide a distorted picture of functional similarity.

In this regard, the results suggest that task-based model stitching is inferior to alignment-based approaches. Specifically, we observe that both SLM and TLM not only exploit shortcuts but also overlook existing functional alignments, leading to degraded performance on clean data. In contrast, alignment-based model stitching retains its ability to perform well on clean data, suggesting that it captures functional similarity more effectively. However, we also observe that even Hint and FuLA are susceptible to shortcut learning, particularly in the absence of functional similarity. For instance, location shortcuts are exploited when stitching deeper layers of ImageNet models to CIFAR models, but not when stitching independently trained models on CIFAR.

Severe discrepancies between shortcut and clean data indicate overfitting to the training distribution, in this case, the shortcuts. Based on this, it is evident that task-based model stitching is prone to overfitting in the presence of covariance shift (e.g. shortcuts). To further relate model stitching with overfitting, we consider a setting of label shift. In particular, we perform cross-task model stitching using a subset of CIFAR-10 [3]. Fig. 9 suggests that all model stitching methods underperform when trained on samples from a limited number of classes. However, we observe a clear trend: TLM and SLM tend to overfit more than Hint and FuLA. In this setting, FuLA emerges as the most effective model stitching approach, maintaining high functional similarity across the full distribution despite being trained on only a subset of it. Notably, in the same-task



Figure 9: Model stitching when DM-initialized and trained on a subset of classes.

stitching setup, FuLA outperformed the others only in the single-class setting. Beyond that, SLM, Hint, and FuLA performed comparably well, while TLM consistently underperformed.

**Takeaway:** Alignment-based model stitching identifies more meaningful notions of functional similarity, as it is less sensitive to the data used for alignment compared to task-based model stitching.
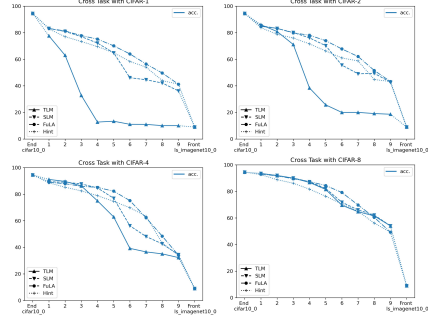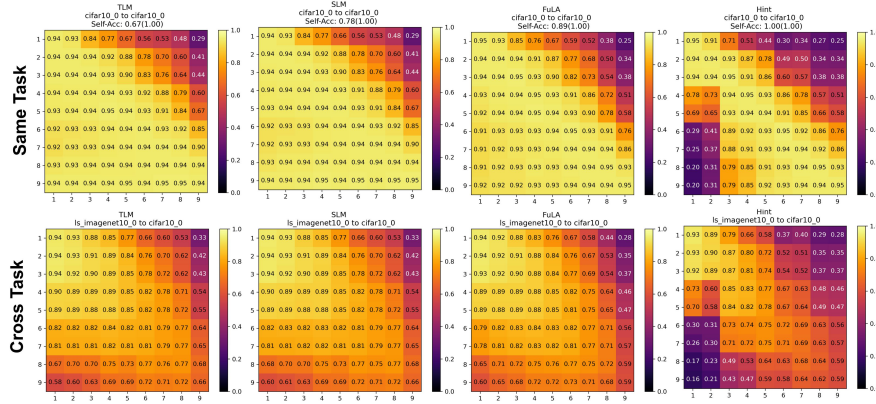
## 3.3 Cross Layer Stitching



Figure 10: Cross-layer stitching between regularly trained models. Rows and columns represent the output layers of the front model and input layers of the end model, respectively. Each cell shows the stitched model's accuracy (i.e., functional similarity). Self-Accuracy indicates how often a layer is most similar to itself, with parentheses denoting non-directional mode performance.

Recently, Balogh and Jelasity [2] questioned the reliability of model stitching under TLM, criticizing its directional nature where high functional similarity is attainable when stitching deep to swallower layers but not vice versa [19]. A consequence to that is high functional similarity between distant layers sometimes exceeds that of architecturally corresponding ones, even within identical networks. In response, DM was proposed as a more appropriate similarity metric.

Previously, we have seen that task-based stitching can be unreliable due to its reliance on output level cues to solve the task, and in that sense directionality could potentially be a byproduct of that. We revisit cross-layer stitching and find that task-agnostic FuLA also behaves in a directional manner. In particular, we have observed a smooth transition as we included deeper functional terms (i.e., through weighting the FuLA objective the Eq. (5)) as shown in Fig. 11. Based on that, we conclude that earlier layers can retain functionally relevant information when guided by deeper layers. In contrast, guidance from deeper layers, when given swallow-level input is not sufficient to establish functional comparability.

---

[3]Classes were added progressively based on their label index, i.e. CIFAR-1 consists of the 'airplane' class.

We observe that model stitching by Hint (i.e., realizing the DM) also exhibits directional behavior (e.g., [5,3] is higher than [3,5] in self-stitching in Fig. 10) which breaks as we move from the diagonal towards the bottom-left. In relation to Fig. 11, this appears to be a form of regularization on functional alignment enforcing non-directional behavior. We argue that directionality can also break trivially by considering both directions to infer similarity. For example in Fig 10 corresponding layers in self-stitching are accurately detected, when incorporating (i.e., averaging) both directions to identify them.



Figure 11: Cross-layer stitching with FuLA when including up to three functional Hints.

On the other hand, the directionality of functional similarity can provide fine-grained insights. For example, when comparing the rows in Fig. 10 under FuLA, it appears that skipping the first 5 in-distribution layers is functionally equivalent to using the penultimate representations of an out-of-distribution task (i.e., these settings achieved 59% and 60% accuracy). Another interesting observation is that functional similarity appears to diffuse gradually across the axons of task specialization (see 2nd in Fig. 10) which is a property missed by Hint, but maintained for the rest of the settings.

**Takeaway:** Model stitching is a directional functional similarity metric even when inferred using task-agnostic functional alignment enabled by FuLA.
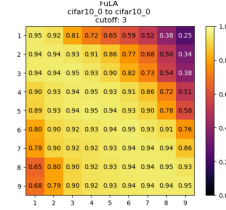
## 4   Related Work

**Model Stitching:** Model stitching was introduced by Lenc and Vedaldi [25], who used a $1 \times 1$ affine transformation to connect intermediate convolutional layers to study invariance and equivariance. This was later extended by [3, 9] to measure functional similarity. Recently, the reliability of model stitching for functional similarity has been questioned, arguing that high similarity scores may arise even among models with different biases [40] or between distant layers [19]. Beyond functional similarity, model-stitching has been used to connect independent network components from the model zoo [37, 38, 43] allowing for networks of controllable performance-efficiency trade-offs. Zero-shot stitching has also been achieved by learning representations in relative space [5, 34].

**Feature distillation:** Feature-based distillation was introduced by Romero et al. [39], where teacher features guide the student via L2 minimization. Notable extensions include attention-based distillation [44] and inter-channel correlation transfer [28], both applied at manually chosen layers. Ji et al. [20] and Chen et al. [6] avoid manual layer linking by learning connections via attention. While FuLA also minimizes L2 feature distance, it does so through frozen layer compositions of the end model, achieving alignment in the functional sense.

**Function consistency:** Function consistency has been used in KD [27, 42] to help students better replicate teacher behavior, improving performance. Yang et al. [43] applied it to group equivalent blocks in the model zoo, later reassembled under resource constraints. In contrast, FuLA leverages function consistency to align models for measuring functional similarity.

## 5   Conclusions

In this work, we proposed FuLA as a novel optimization condition for model stitching as a functional similarity metric. We show that stitching models using output-level cues (i.e., SLM and TLM) can fabricate alignment by overfitting to the training task. Experiments under adversarial training, distribution shifts, and cross-task stitching demonstrate that FuLA is significantly less susceptible to pitfalls of task-based stitching. Moreover, it retrieves non-trivial alignment transformations missed by stitch-level alignment, making it a more reliable approach for evaluating functional similarity.

**Limitations:** While FuLA shows promise as a more reliable metric for functional similarity compared to prior model stitching approaches, definitive conclusions remain out of reach due to the lack to ground truth. Furthermore, our study is limited to image classification, exploring FuLA's relation to task-based stitching in other modalities and tasks is a promising direction for future work.

# References

[1] András Balogh and Márk Jelasity. On the functional similarity of robust and non-robust neural representations. In *International Conference on Machine Learning*, pages 1614–1635. PMLR, 2023.

[2] András Balogh and Márk Jelasity. How not to stitch representations to measure similarity: Task loss matching versus direct matching. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 39, pages 15472–15480, 2025.

[3] Yamini Bansal, Preetum Nakkiran, and Boaz Barak. Revisiting model stitching to compare neural representations. *Advances in neural information processing systems*, 34:225–236, 2021.

[4] Yoshua Bengio, Aaron Courville, and Pascal Vincent. Representation learning: A review and new perspectives. arxiv 2012. *arXiv preprint arXiv:1206.5538*, 2012.

[5] Irene Cannistraci, Luca Moschella, Marco Fumero, Valentino Maiorca, and Emanuele Rodolà. From bricks to bridges: Product of invariances to enhance latent space communication. *arXiv preprint arXiv:2310.01211*, 2023.

[6] Defang Chen, Jian-Ping Mei, Yuan Zhang, Can Wang, Zhe Wang, Yan Feng, and Chun Chen. Cross-layer distillation with semantic calibration. In *Proceedings of the AAAI conference on artificial intelligence*, volume 35, pages 7028–7036, 2021.

[7] Rhea Chowers and Yair Weiss. What do cnns learn in the first layer and why? a linear systems perspective. In *International Conference on Machine Learning*, pages 6115–6139. PMLR, 2023.

[8] Francesco Croce and Matthias Hein. Reliable evaluation of adversarial robustness with an ensemble of diverse parameter-free attacks. In *International conference on machine learning*, pages 2206–2216. PMLR, 2020.

[9] Adrián Csiszárik, Péter Kőrösi-Szabó, Akos Matszangosz, Gergely Papp, and Dániel Varga. Similarity and matching of neural network representations. *Advances in Neural Information Processing Systems*, 34:5656–5668, 2021.

[10] MohammadReza Davari, Stefan Horoi, Amine Natik, Guillaume Lajoie, Guy Wolf, and Eugene Belilovsky. Reliability of cka as a similarity measure in deep learning. *arXiv preprint arXiv:2210.16156*, 2022.

[11] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009.

[12] Frances Ding, Jean-Stanislas Denain, and Jacob Steinhardt. Grounding representation similarity through statistical testing. *Advances in Neural Information Processing Systems*, 34:1556–1568, 2021.

[13] Diego Doimo, Aldo Glielmo, Alessio Ansuini, and Alessandro Laio. Hierarchical nucleation in deep neural networks. *Advances in Neural Information Processing Systems*, 33:7526–7536, 2020.

[14] Andre Esteva, Brett Kuprel, Roberto A Novoa, Justin Ko, Susan M Swetter, Helen M Blau, and Sebastian Thrun. Dermatologist-level classification of skin cancer with deep neural networks. *nature*, 542(7639):115–118, 2017.

[15] Sergio Garrido-Jurado, Rafael Muñoz-Salinas, Francisco José Madrid-Cuevas, and Manuel Jesús Marín-Jiménez. Automatic generation and detection of highly reliable fiducial markers under occlusion. *Pattern Recognition*, 47(6):2280–2292, 2014.

[16] Robert Geirhos, Jörn-Henrik Jacobsen, Claudio Michaelis, Richard Zemel, Wieland Brendel, Matthias Bethge, and Felix A Wichmann. Shortcut learning in deep neural networks. *Nature Machine Intelligence*, 2(11):665–673, 2020.

[17] Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. *arXiv preprint arXiv:1412.6572*, 2014.

[18] 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 recognition*, pages 770–778, 2016.

[19] Adriano Hernandez, Rumen Dangovski, Peter Y. Lu, and Marin Soljacic. Model stitching: Looking for functional similarity between representations. In *SVRHM 2022 Workshop @ NeurIPS*, 2022. URL https://openreview.net/forum?id=7bHLCO5FQdB.

[20] Mingi Ji, Byeongho Heo, and Sungrae Park. Show, attend and distill: Knowledge distillation via attention-based feature matching. In *Proceedings of the AAAI conference on artificial intelligence*, volume 35, pages 7945–7952, 2021.

[21] Haydn T Jones, Jacob M Springer, Garrett T Kenyon, and Juston S Moore. If you've trained one you've trained them all: inter-architecture similarity increases with robustness. In *Uncertainty in Artificial Intelligence*, pages 928–937. PMLR, 2022.

[22] Diederik P Kingma. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.

[23] Simon Kornblith, Mohammad Norouzi, Honglak Lee, and Geoffrey Hinton. Similarity of neural network representations revisited. In *International conference on machine learning*, pages 3519–3529. PMLR, 2019.

[24] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. *Communications of the ACM*, 60(6):84–90, 2017.

[25] Karel Lenc and Andrea Vedaldi. Understanding image representations by measuring their equivariance and equivalence. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 991–999, 2015.

[26] Tao Li, Yingwen Wu, Sizhe Chen, Kun Fang, and Xiaolin Huang. Subspace adversarial training. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13409–13418, 2022.

[27] Dongyang Liu, Meina Kan, Shiguang Shan, and Xilin Chen. Function-consistent feature distillation. *arXiv preprint arXiv:2304.11832*, 2023.

[28] Li Liu, Qingle Huang, Sihao Lin, Hongwei Xie, Bing Wang, Xiaojun Chang, and Xiaodan Liang. Exploring inter-channel correlation for diversity-preserved knowledge distillation. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 8271–8280, 2021.

[29] Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. Towards deep learning models resistant to adversarial attacks. *arXiv preprint arXiv:1706.06083*, 2017.

[30] Wojciech Masarczyk, Mateusz Ostaszewski, Ehsan Imani, Razvan Pascanu, Piotr Miłoś, and Tomasz Trzcinski. The tunnel effect: Building data representations in deep neural networks. *Advances in Neural Information Processing Systems*, 36:76772–76805, 2023.

[31] Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. Efficient estimation of word representations in vector space. *arXiv preprint arXiv:1301.3781*, 2013.

[32] Seyed Iman Mirzadeh, Mehrdad Farajtabar, Dilan Gorur, Razvan Pascanu, and Hassan Ghasemzadeh. Linear mode connectivity in multitask and continual learning. *arXiv preprint arXiv:2010.04495*, 2020.

[33] Ari Morcos, Maithra Raghu, and Samy Bengio. Insights on representational similarity in neural networks with canonical correlation. *Advances in neural information processing systems*, 31, 2018.

[34] Luca Moschella, Valentino Maiorca, Marco Fumero, Antonio Norelli, Francesco Locatello, and Emanuele Rodolà. Relative representations enable zero-shot latent space communication. *arXiv preprint arXiv:2209.15430*, 2022.

[35] Vedant Nanda, Till Speicher, Camila Kolling, John P Dickerson, Krishna Gummadi, and Adrian Weller. Measuring representational robustness of neural networks through shared invariances. In *International Conference on Machine Learning*, pages 16368–16382. PMLR, 2022.

[36] Thao Nguyen, Maithra Raghu, and Simon Kornblith. Do wide and deep networks learn the same things? uncovering how neural network representations vary with width and depth. *arXiv preprint arXiv:2010.15327*, 2020.

[37] Zizheng Pan, Jianfei Cai, and Bohan Zhuang. Stitchable neural networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16102–16112, 2023.

[38] Zizheng Pan, Jing Liu, Haoyu He, Jianfei Cai, and Bohan Zhuang. Stitched vits are flexible vision backbones. In *European Conference on Computer Vision*, pages 258–274. Springer, 2024.

[39] Adriana Romero, Nicolas Ballas, Samira Ebrahimi Kahou, Antoine Chassang, Carlo Gatta, and Yoshua Bengio. Fitnets: Hints for thin deep nets. *arXiv preprint arXiv:1412.6550*, 2014.

[40] Damian Smith and Antonia Marcu. Comparing apples and oranges: is stitching similarity a load of spheres? In *NeurIPS 2024 Workshop on Scientific Methods for Understanding Deep Learning*, 2024.

[41] Dimitris Tsipras, Shibani Santurkar, Logan Engstrom, Alexander Turner, and Aleksander Madry. Robustness may be at odds with accuracy. *arXiv preprint arXiv:1805.12152*, 2018.

[42] Jing Yang, Brais Martinez, Adrian Bulat, and Georgios Tzimiropoulos. Knowledge distillation via softmax regression representation learning. In *International Conference on Learning Representations*, 2021.

[43] Xingyi Yang, Daquan Zhou, Songhua Liu, Jingwen Ye, and Xinchao Wang. Deep model reassembly. *Advances in neural information processing systems*, 35:25739–25753, 2022.

[44] Sergey Zagoruyko and Nikos Komodakis. Paying more attention to attention: Improving the performance of convolutional neural networks via attention transfer. *arXiv preprint arXiv:1612.03928*, 2016.

[45] Liu Ziyin, Isaac Chuang, Tomer Galanti, and Tomaso Poggio. Formation of representations in neural networks. *arXiv preprint arXiv:2410.03006*, 2024.

# A Experimental Setup

In this section, we cover the implementation details needed for reproducing our results.

## A.1 Baseline Training

In this study, we experimented with CNN-based architectures, namely ResNet18, ResNet34, and a modified version of ResNet18 in which the residual connections were removed, effectively rendering it a plain convolutional neural network (CNN). To train the end and front models we used the CIFAR-10 and low resolution version of ImageNet where we randomly sampled 10 classes, we refer to that as LS-ImageNet-10.

**Non-robust models:** The non-robust models were trained for 200 epochs using the SGD optimizer with momentum. The initial learning rate was set to 0.1, with a batch size of 128. During training, we applied random horizontal flips and random cropping with 4-pixel padding. A weight decay of $1 \times 10^{-4}$ was used, and the learning rate was reduced by a factor of 0.1 at every third of the total training epochs.

**Robust models:** The robust models were trained using the same hyperparameter settings as the non-robust models, with adversarial training performed according to [29], using an equal mix of clean and adversarial examples ($\alpha = 0.5$). The adversarial examples were generated using 10 steps, with a step size of $2/255$ and $\epsilon$ set to $8/255$. We follow [1] and measure adversarial robustness was measured using the "rand" configuration of AutoAttack [8].

The performance of each model is shown in Tab 1.

Table 1: Performance for the baseline models used as front and end models during stitching.

| Model ID | Dataset | Architecture | Accuracy | Robust Accuracy |
|---|---|---|---|---|
| cifar10_0 | CIFAR-10 | ResNet18 (w/o residuals) | 94.21% | 0% |
| | | ResNet18 | 94.51% | 0% |
| | | ResNet34 | 94.14% | 0% |
| cifar10_1 | CIFAR-10 | ResNet18 (w/o residuals) | 94.21% | 0% |
| | | ResNet18 | 94.03% | 0% |
| | | ResNet34 | 94.91% | 0% |
| ls_imagenet10_0 | LS-ImageNet-10 | ResNet18 (w/o residuals) | 76.12% | 0% |
| | | ResNet18 | 74.8% | 0% |
| | | ResNet34 | 74.93% | 0% |
| cifar10_0 robust | CIFAR-10 w/ AT | ResNet18 (w/o residuals) | 84.86% | 40.59% |
| | | ResNet18 | 85.52% | 40.73% |
| | | ResNet34 | 86.69% | 42.89% |
| cifar10_1 robust | CIFAR-10 w/ AT | ResNet18 (w/o residuals) | 85.17% | 40.99% |
| | | ResNet18 | 84.96% | 40.65% |
| | | ResNet34 | 86.22% | 42.86% |
| ls_imagenet10_0 robust | LS-ImageNet-10 w/ AT | ResNet18 (w/o residuals) | 66.2% | 21.8% |
| | | ResNet18 | 63.23% | 22.6% |
| | | ResNet34 | 63.79% | 22.2% |

## A.2 Stitching Training

We follow [1, 9] and use $1 \times 1$ convolutional layers with bias as a stitch layer to connect the front and end models. During training we keep everything frozen apart from the alignment transformation and the batch normalization statistics of the front and end models. The affine transformation was

trained for 30 epochs using Adam [22] optimizer with a learning rate of $0.001$ and weight decay of $1 \times 10^{-4}$. During training, we use the same augmentation strategy used to train the base models. We perform stitching at the first convolutional layer and at the first and last (residual) block of each layer leading to 9 stitching locations for all architectures considered. At all times, the stitching layer was initialized by DM using 100 samples for the training set.

**Model Stitching under Adversarial Training:** When performing stitching under adversarial training, we consider three $\alpha$ configurations, following Balogh and Jelasity [1]: $\alpha = 0$ (stitching using only clean samples), $\alpha = 0.5$ (stitching using an equal mix of clean and adversarial examples), and $\alpha = 1$ (stitching using only adversarial examples). The adversarial examples were generated for the stitched model using the same hyperparameters as those used to train the robust models.

**Model Stitching under Shortcut:** We consider various shortcut settings, namely: (i) patterns that correlate with the class at random locations, (ii) random patterns (uniform noise) at fixed locations correlated with the class, and (iii) a combination of both. For the patterns, we use either $4 \times 4$ squares of fixed values depending on the class—referred to as pixels—or $4 \times 4$ ArUco markers [15]. An example of these shortcut configurations is shown in Fig. 12.



Figure 12: Visualizing shortcut configurations.

**Cross-Layer Stitching:** During cross-layer stitching, we train the affine transformation using the same training configuration as used for the corresponding layer stitching. We performed stitching for all $9 \times 9$ layer combinations. When stitching layers that produce feature maps of different spatial resolutions, we follow Balogh and Jelasity [2] and use bilinear interpolation to match the resolution of the front feature maps to that of the end feature maps. To convert directional functional similarity into a non-directional measure, we averaged the functional similarities in both stitching directions. For example, the similarity between the $8^{\text{th}}$ and $1^{\text{st}}$ layers is computed as the average of the functional similarities when stitching from the $1^{\text{st}}$ layer to the $8^{\text{th}}$, and from the $8^{\text{th}}$ to the $1^{\text{st}}$.

# B Full Results on Additional Architectures

In this section, we provide experimental results on additional architectures and shortcut learning settings.

## B.1 Model Stitching under Adversarial Training

When it comes to model stitching under adversarial training, we observe that identical conclusions can be drawn across all CNN architectures: both SLM and TLM stitching settings can fabricate alignment sufficient to solve the task. For example, task-based stitching can report high functional similarity with respect to adversarial accuracy even when stitching between non-robust models, while FuLA and Hint can achieve non-zero adversarial robustness only when at least one of the stitched models is robust. Additionally, in almost all cases, FuLA outperforms Hint in stitching performance. As discussed in the main paper, cases where FuLA underperforms compared to Hint may indicate functional misalignment at deeper layers, which Hint may fail to detect.
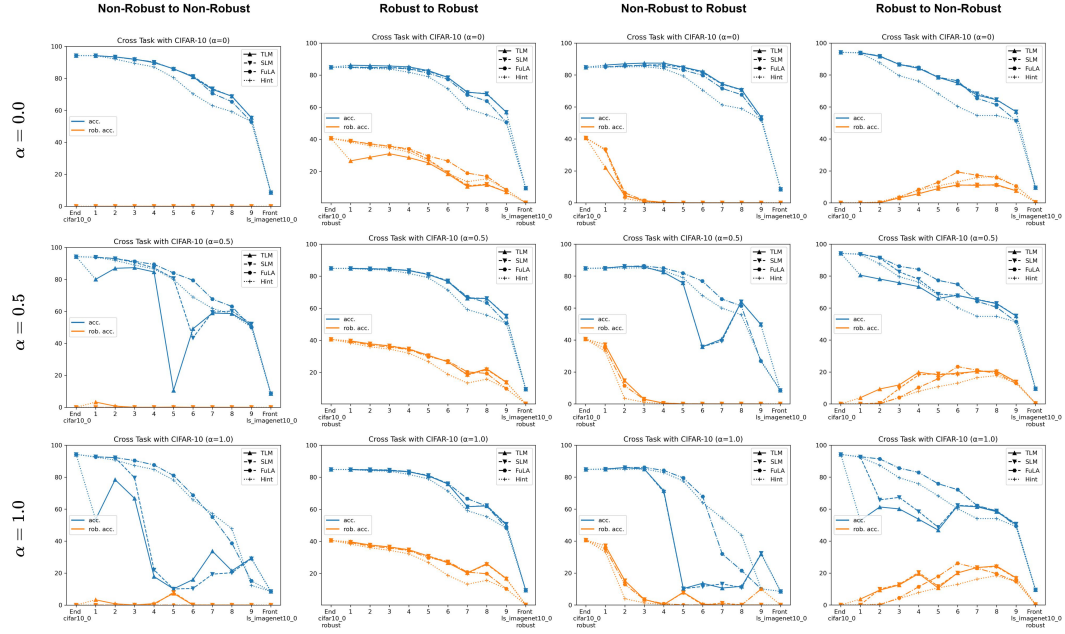
Figure 13: Cross-task stitching under AT – ResNet18 (w/o residuals).
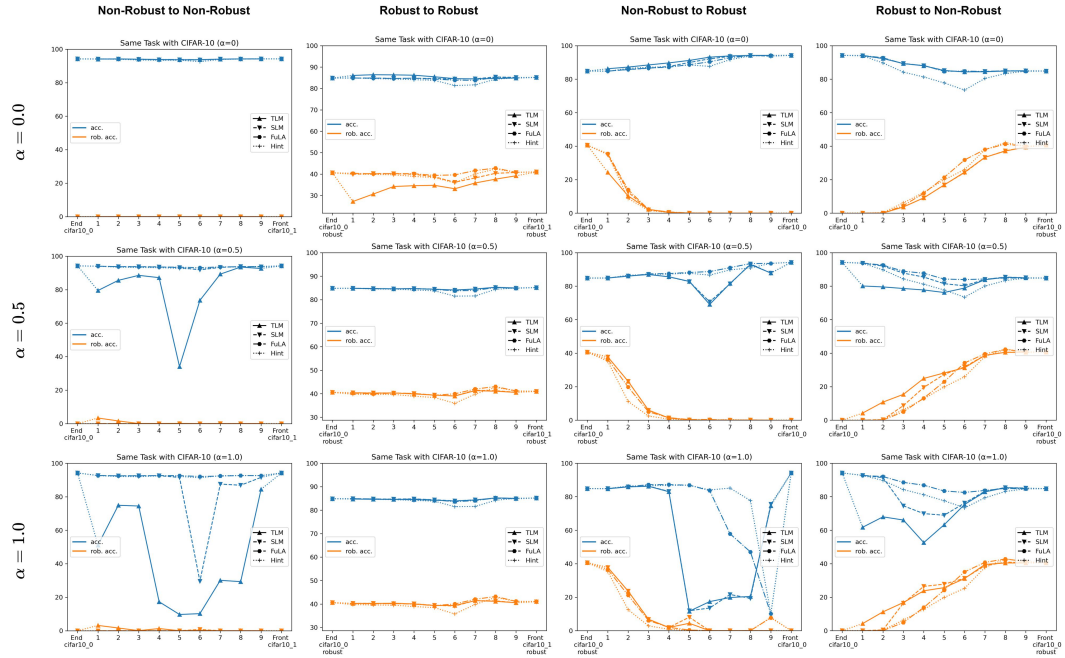


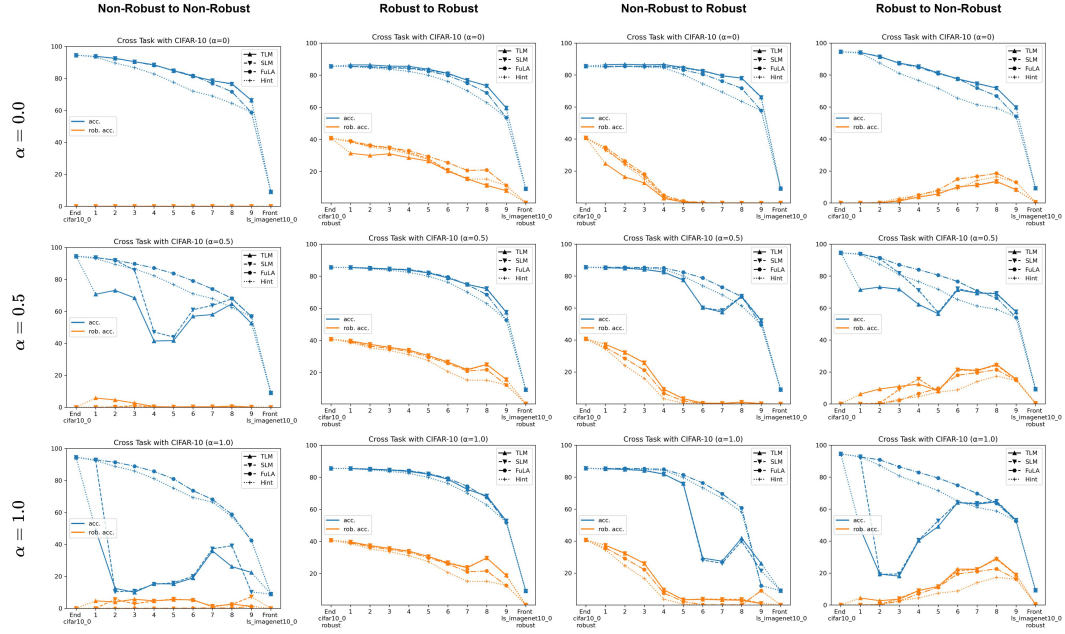Figure 14: Same-task stitching under AT – ResNet18 (w/o residuals).

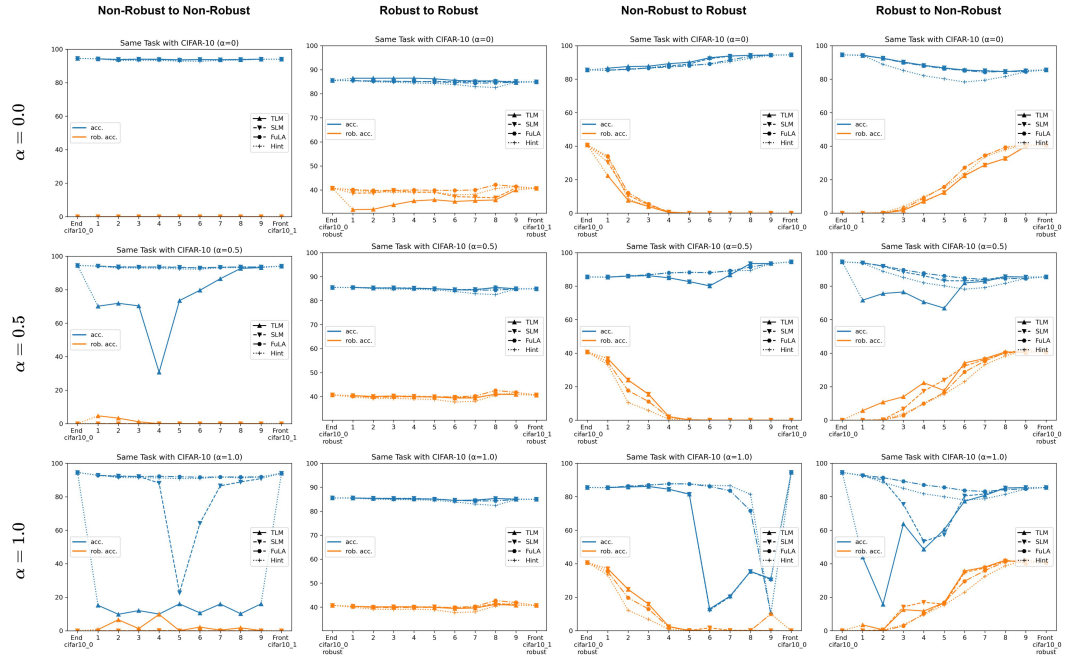Figure 15: Cross-task stitching under AT – ResNet18.
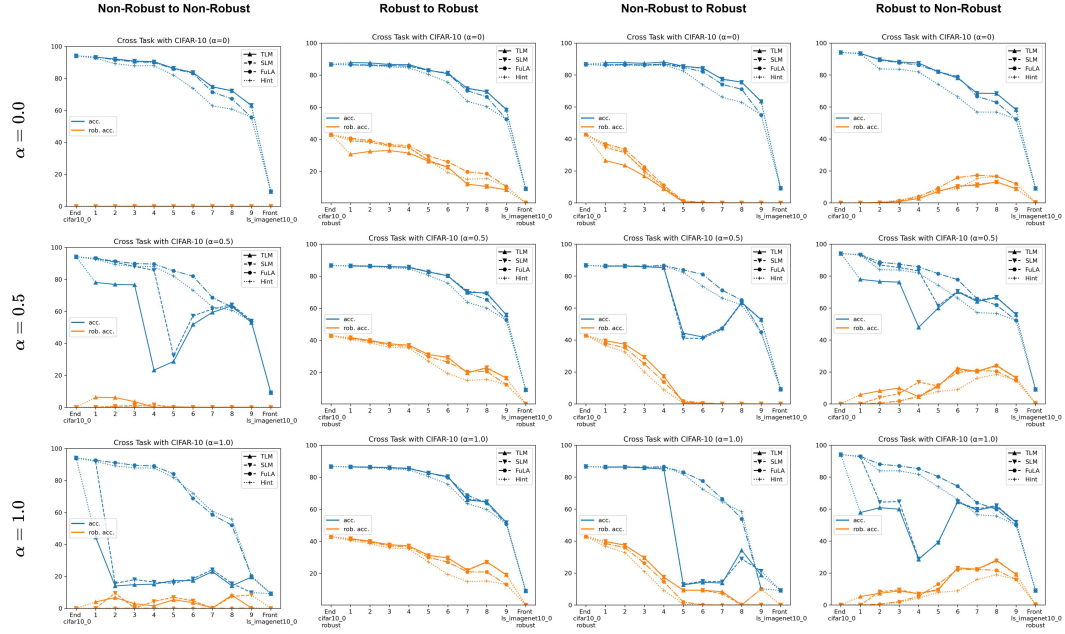


Figure 16: Same-task stitching under AT – ResNet18.

16
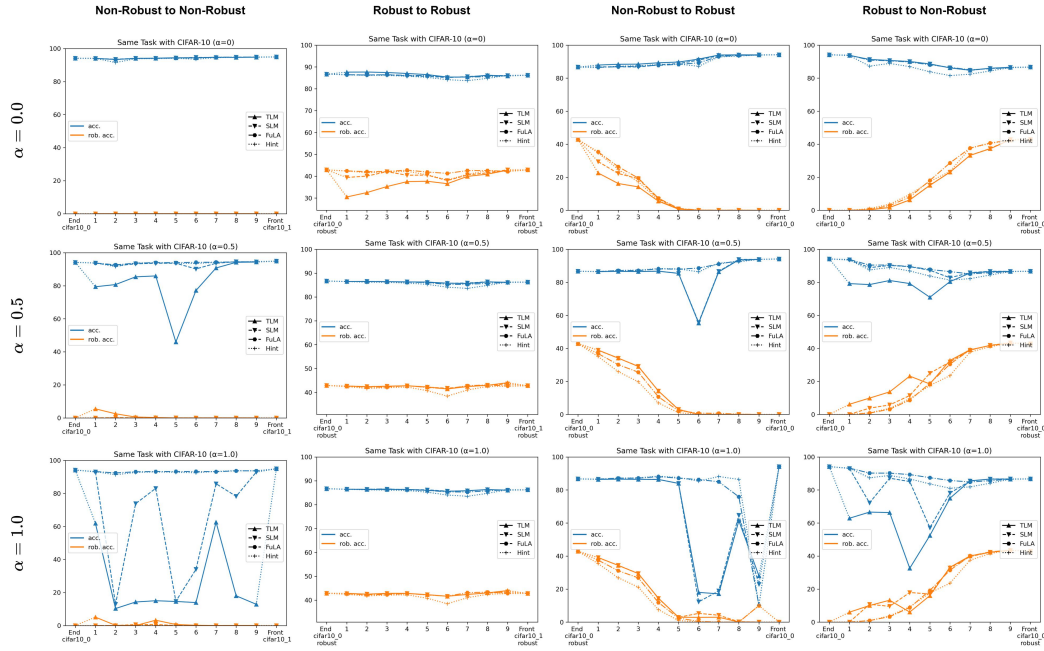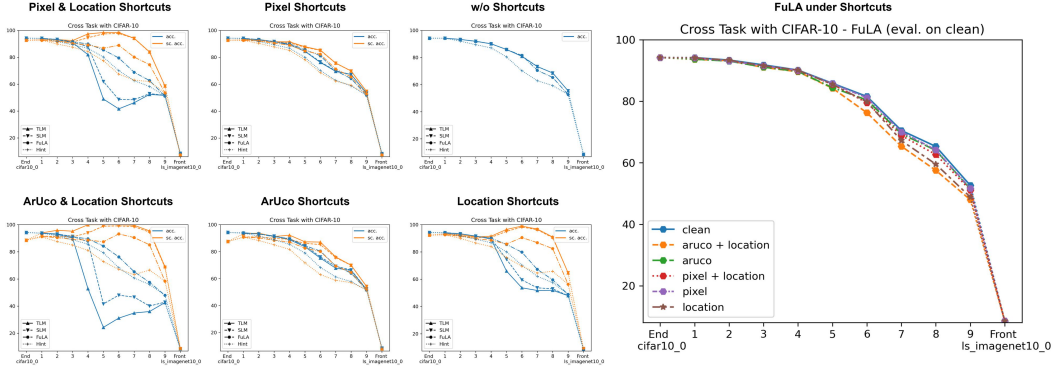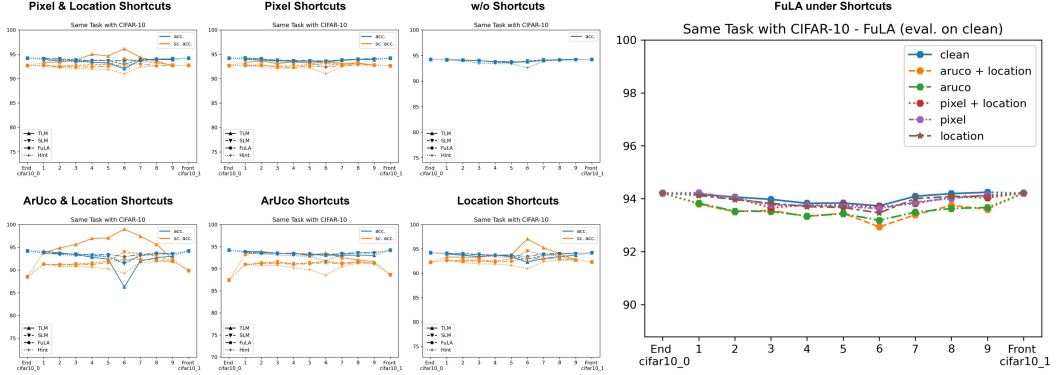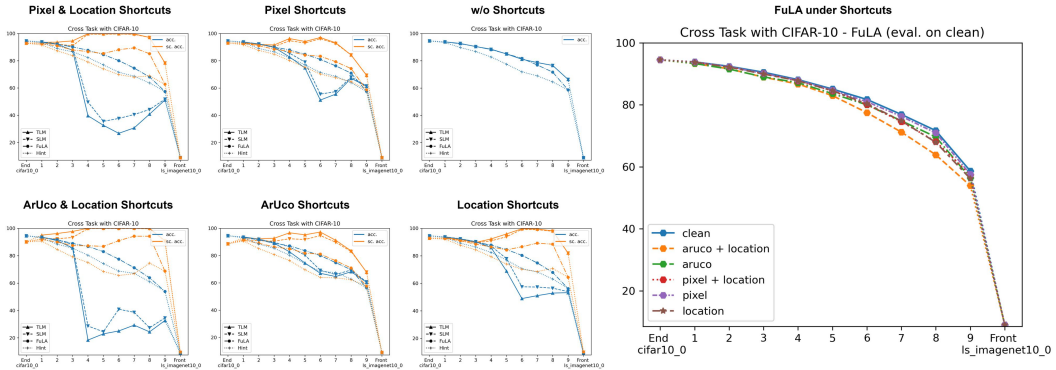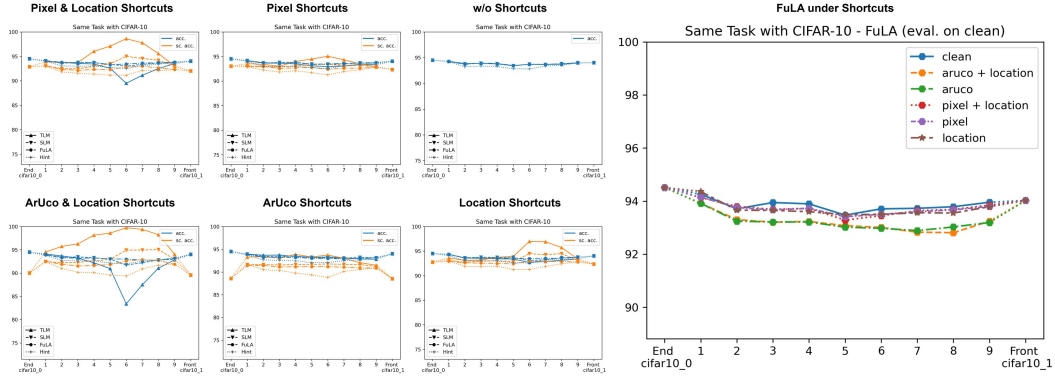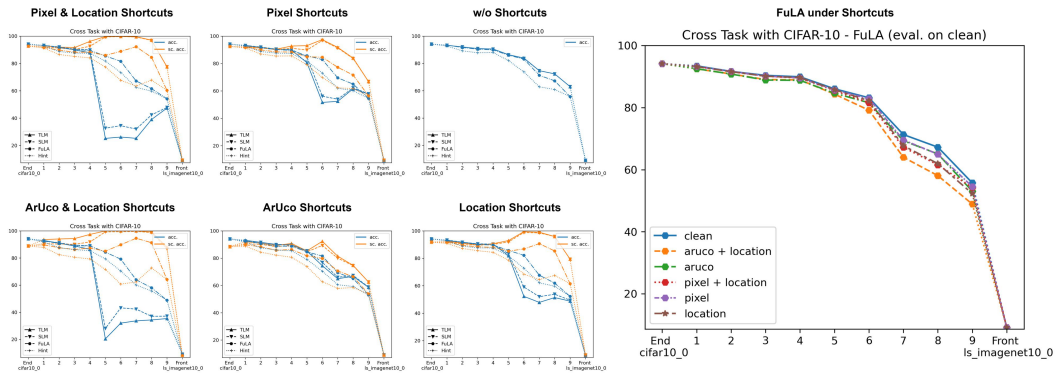
Figure 17: Cross-task stitching under AT – ResNet34.



Figure 18: Same-task stitching under AT – ResNet34.

## B.2 Model Stitching under Shortcuts

We observe that all architectures considered exhibit similar behavior across the different shortcut settings. In particular, task-based stitching tends to overfit to the shortcuts—achieving high performance on them but reduced performance on clean data—indicating that both SLM and TLM opt for alignments that maximize task performance, which can be problematic for evaluating functional similarity. Once again, FuLA outperforms Hint, while performance on clean data remains relatively stable across the different shortcut settings. This suggests that FuLA captures a more meaningful notion of functional alignment.



Figure 19: Cross-task stitching under shortcuts – ResNet18 (w/o residuals).



Figure 20: Same-task stitching under shortcuts – ResNet18 (w/o residuals).



Figure 21: Cross-task stitching under shortcuts – ResNet18.

Figure 22: Same-task stitching under shortcuts – ResNet18.



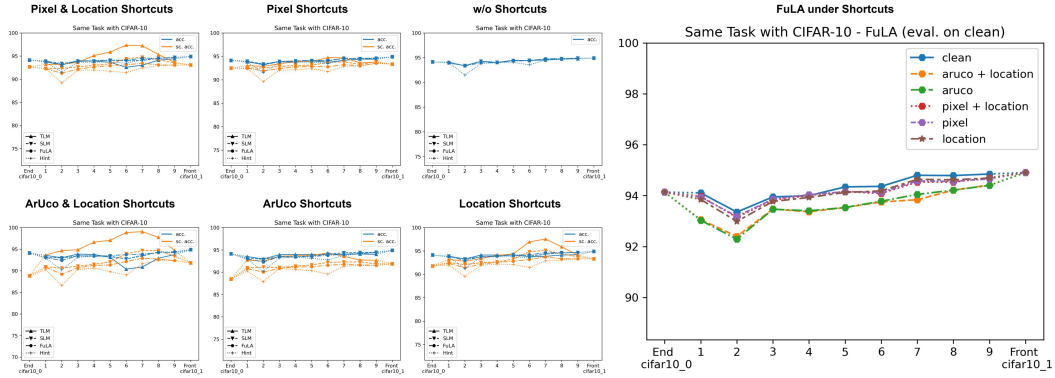Figure 23: Cross-task stitching under shortcuts – ResNet34.



Figure 24: Same-task stitching under shortcuts – ResNet34.

19

## B.3 Model Stitching on a Subset of Classes

When training on a subset of classes, FuLA demonstrates greater robustness, where alignment on the subset leads to improved stitching performance across all classes. This effect is more pronounced in the cross-task setting, but it is also observable in the same-task setting.



Figure 25: Stitching on a subset of classes – ResNet18 (w/o residuals).
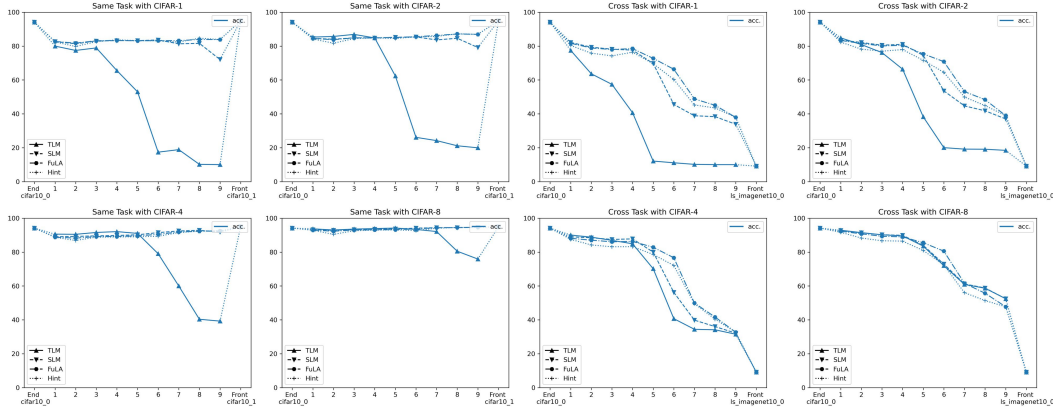


Figure 26: Stitching on a subset of classes – ResNet18.



Figure 27: Stitching on a subset of classes – ResNet34.

## B.4 Cross Layer Stitching

Cross-layer stitching reveals that the directionality of functional similarity is not induced by the task, as we also observe it when stitching under FuLA, which is task-agnostic. When applicable, enforcing non-directionality at a post-hoc stage appears to improve the corresponding layer detection performance across all settings. As expected, the benefit of converting Hint into a non-directional metric is smaller, since Hint is already regularized to be less directional during training.
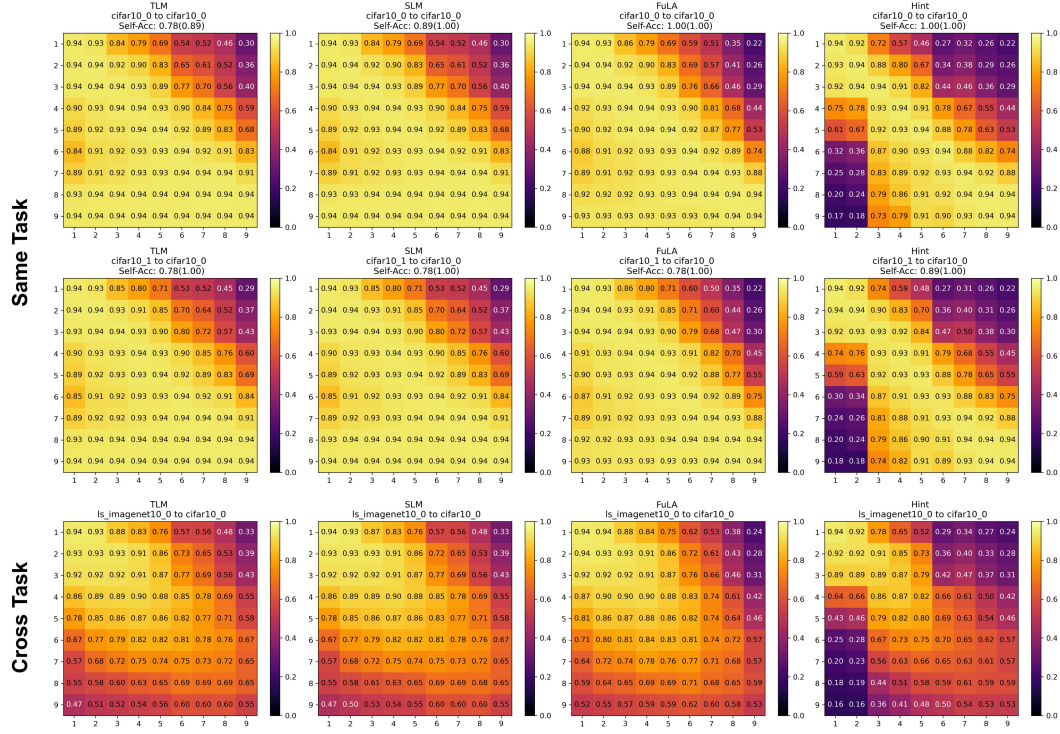


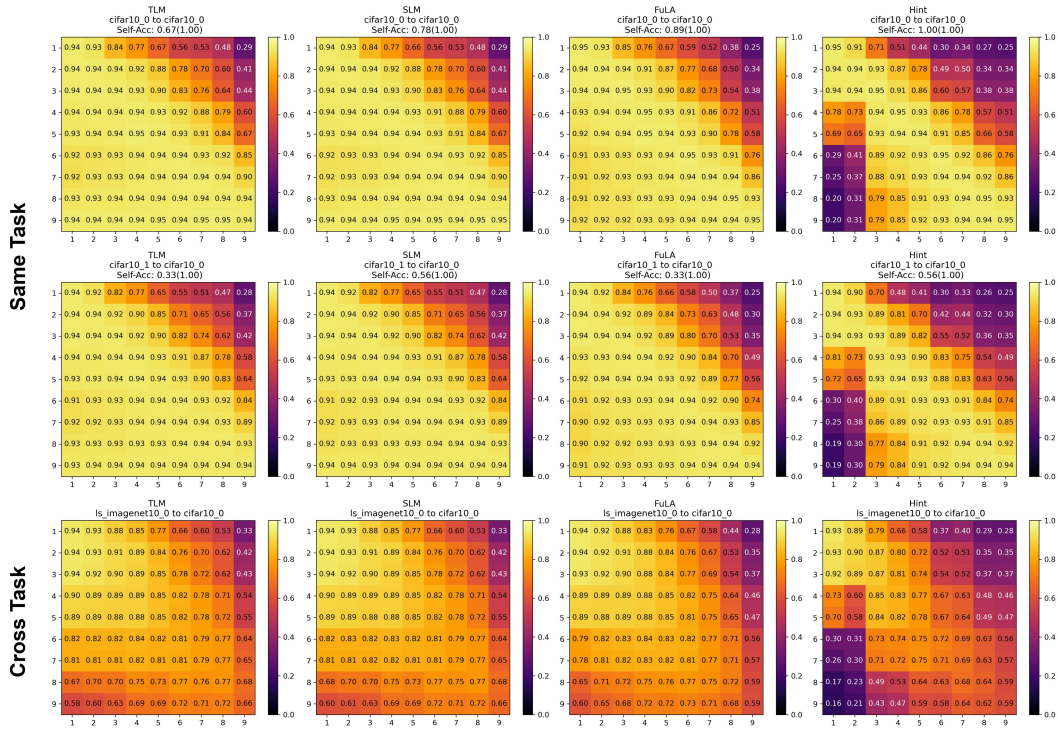Figure 28: Cross-layer stitching – ResNet18 (w/o residuals).
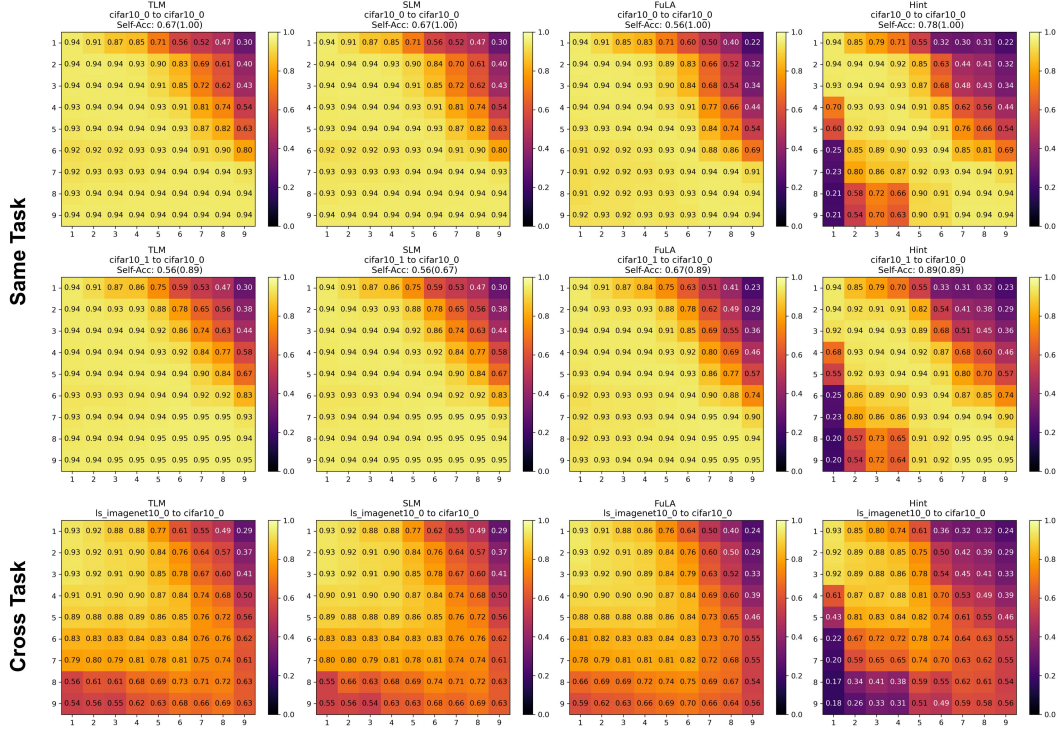
Figure 29: Cross-layer stitching – ResNet18.



Figure 30: Cross-layer stitching – ResNet34.