

000
001
002
003
004
005
006
007
008
009
010
011
012
013
014
015
016
017
018
019
020
021
022
023
024
025
026
027
028
029
030
031
032
033
034
035
036
037
038
039
040
041
042
043
044
045
046
047
048
049
050
051
052
053
054

Convergent World Representations and Divergent Tasks

Anonymous Authors¹

Abstract

While neural representations are central to modern deep learning, the conditions governing their geometry and their roles in downstream adaptability remain poorly understood. We develop a framework clearly separating the underlying world, the data generation process and the resulting model representations to study these questions in a controlled setup. 5,075 city coordinates define the world and 7 geometric tasks generate the training data for autoregressive training. We find that different tasks give rise to qualitatively and quantitatively distinct world representation geometries. However, multi-task training drives convergence of world representations: models trained on non-overlapping tasks develop aligned geometric representations, providing controlled evidence for the Multitask Scaling Hypothesis of the Platonic Representation Hypothesis. To study adaptation, we pretrain models on all tasks, then test whether *new* entities (cities) can be consistently integrated into the representation space via fine-tuning. Surprisingly, we find that despite multi-task pretraining, some tasks, which we call *divergent*, actively harm the representational integration of new entities and harm generalization. Our results show that training on multiple relational tasks reliably produces convergent world representations, but lurking divergent tasks can catastrophically harm new entity integration via fine-tuning.

1. Introduction

The nature of representations and mechanisms learned by deep neural networks, or in fact any intelligent system, and their relation to generalization is a central topic in deep learning research (Hubel & Wiesel, 1962; Rosenblatt, 1958;

Fukushima, 1980; Rumelhart et al., 1986). Recent work has demonstrated that neural networks trained on vast amounts of data can capture diverse, disentangled and sometimes interpretable aspects of their training data, or even of the world underlying the data (Bengio et al., 2014). These rich representations are generally thought to underlie the generalization and adaptability of neural networks to unseen, out-of-distribution scenarios.

Recent work on internal representations of language models has provided evidence that neural networks can develop structured representations of the underlying data rather than merely memorizing surface patterns (Li et al., 2022; Gurnee & Tegmark, 2023; Nanda et al., 2023b).

However, major open questions remain. When interpretable representations are discovered in neural networks, it is often unclear whether their emergence is surprising or inevitable, what geometry they will take and how they support generalization. Even less understood is how these representations adjust during fine-tuning and downstream adaptation.

Answering these questions is difficult in real-world settings, where the key factors, the world, the data and the model, are entangled and costly to vary independently. In this work, we develop a synthetic framework where these factors can be precisely controlled and systematically studied.

This work. To study these questions, we decouple the underlying *world* from the *data generation process* to control them independently. Concretely, we adopt the coordinates of real-world cities as our “world,” a ready-made complex structure with ground-truth geometry, and define 7 geometric tasks on top of it. We train autoregressive Transformers on this data and study how world representations form and vary across tasks, surfacing preliminary evidence for the Platonic Representation Hypothesis (PRH) (Huh et al., 2024). Crucially, this setup allows us to define consistent updates to the world (adding new cities) that produce predictable changes in the data, letting us test whether models can absorb such updates via fine-tuning. Our contributions are as follows:

- **A Framework Decoupling World, Data and Model.** (**Sec. 3**) We separate the underlying world (city coordinates) from the data generation process (7 geometric tasks), enabling systematic study of how different tasks

¹Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

Preliminary work. Under review by the International Conference on Machine Learning (ICML). Do not distribute.

055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083
084
085
086
087
088
089
090
091
092
093
094
095
096
097
098
099
100
101
102
103
104
105
106
107
108
109

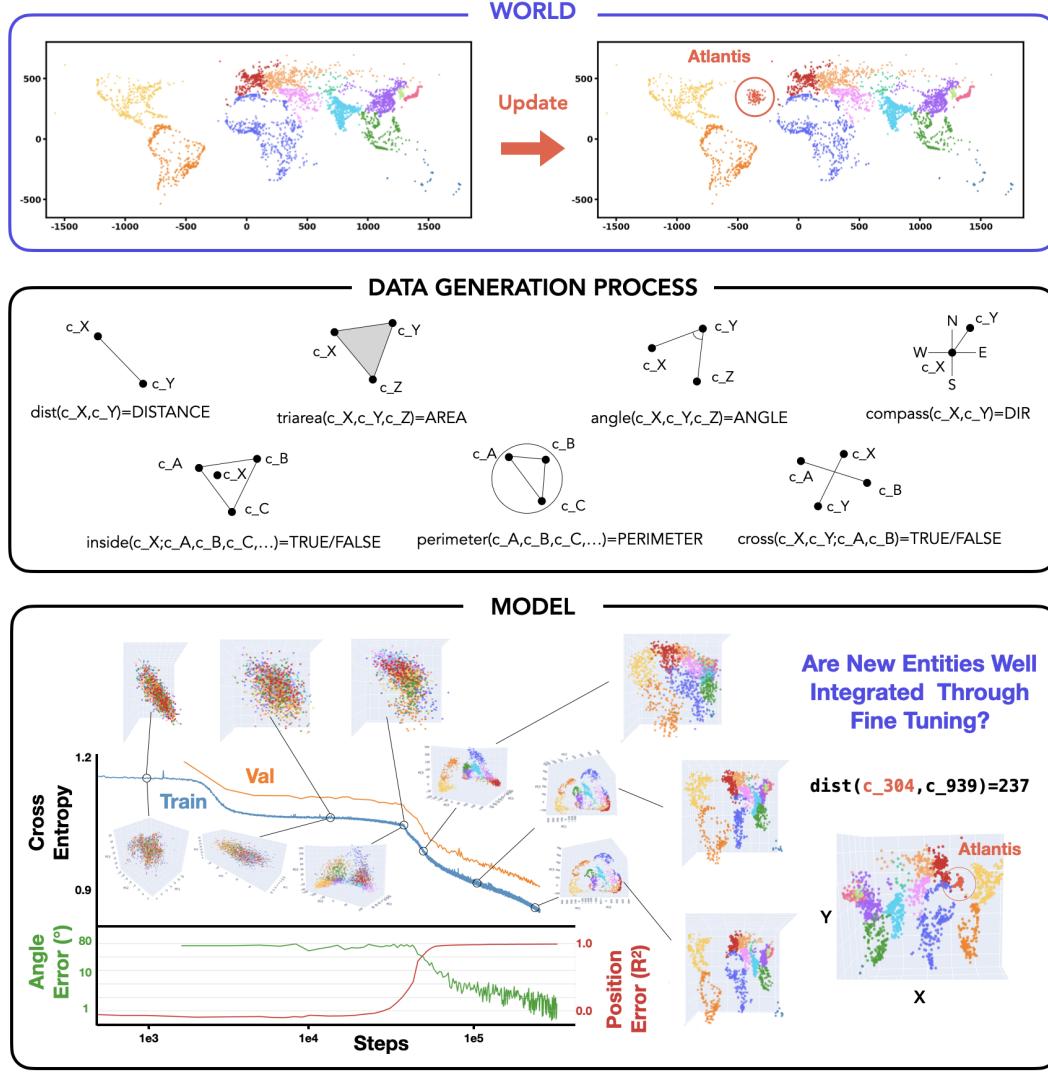


Figure 1. Overview of the World-Data-Model framework. **Top:** The world consists of 5,075 real city coordinates; we test adaptation by adding 100 synthetic Atlantis cities (App. B.1). **Middle:** Seven geometric tasks generate training data from city coordinates (App. B.2). **Bottom:** Training dynamics of one model, showing loss curves, linear probing accuracy for coordinate reconstruction and visualizations of internal representations (PCA and linear probe projections) at different training stages. See App. Fig. 9 for all training curves.

shape representations of the same world. The world provides ground-truth coordinates for directly assessing representation quality via probing. This setup also allows defining consistent world updates (adding synthetic Atlantis cities) to test whether models can adapt their representations accordingly.

- **Task-Dependent Geometry and Multi-Task Convergence. (Sec. 4)** We show that different tasks operating on the same world produce highly variable representational geometries across tasks and seeds. However, multi-task training drives convergence: models trained on multiple tasks show higher representational alignment, even when

they share no common tasks. This provides partial evidence for the Multitask Scaling Hypothesis, one proposed mechanism for the Platonic Representation Hypothesis.

- **Divergent Tasks Harm Fine-Tuning of New Entities Despite Multi-Task Pretraining. (Sec. 5)** We test whether models can integrate new entities (Atlantis cities) via fine-tuning. We find that single-task representational similarity (CKA) partially predicts cross-task generalization. In a multi-task fine-tuning setting, we find surprising “divergent” tasks which hinder integration of new entities into the learned manifold, actively harming generalization.

110 2. Related Work

111 **Internal Representations.** Recent work has revealed that
 112 language models develop structured world models encoding
 113 geographic, temporal and relational information (Li et al.,
 114 2022; Gurnee & Tegmark, 2023; Nanda et al., 2023b; Marks
 115 & Tegmark, 2024). Furthermore, PRH posits that diverse
 116 models converge toward similar representational structures
 117 (Huh et al., 2024), though recent work questions this optimis-
 118 mistic (Kumar et al., 2025). In this work, we study factors
 119 controlling representation formation and how networks inte-
 120 grate new entities via fine-tuning.
 121

122 **Fine-tuning.** The pretraining-finetuning paradigm has be-
 123 come central to modern deep learning. Despite widespread
 124 success, fine-tuning exhibits poorly understood behaviors
 125 such as the reversal curse (Berglund et al., 2024) or emer-
 126 gent misalignment (Betley et al., 2025). On this background,
 127 careful studies of fine-tuning and other low-compute adapta-
 128 tion methods have raised pessimism about whether models
 129 can learn fundamentally new abilities, suggesting they may
 130 merely form “thin wrappers” around pretrained representa-
 131 tions (Jain et al., 2023; Ward et al., 2025; Yue et al., 2025;
 132 Qin et al., 2025). Our work examines this question in a
 133 controlled setup where ground-truth world structure enables
 134 precise measurement of representation adaptation.
 135

136 **Multi-task Learning.** Multi-task learning improves gener-
 137 alization through shared representations (Caruana, 1997);
 138 in some sense, modern foundation models represent an ex-
 139 treme form of multi-task training. Large-scale multi-task
 140 pretraining typically assumes rich representations emerge
 141 from data diversity (Aghajanyan et al., 2021), but the precise
 142 mechanisms remain underexplored. Recent work studies
 143 task diversity in controlled settings (Michaud et al., 2023;
 144 Zhang et al., 2025), though most focus on aggregate behav-
 145 iors rather than characterizing tasks. Here, we define tasks
 146 as geometric functions over a shared world to investigate
 147 how task structure shapes representations.
 148

149 **Synthetic Data.** The cost and complexity of foundation
 150 models has motivated synthetic approaches for controlled
 151 study of in-context learning (Xie et al., 2021; Chan et al.,
 152 2022; Reddy, 2023; Raventós et al., 2023; Park et al., 2024b;
 153 Wurgafit et al., 2025), compositional generalization (Okawa
 154 et al., 2024; Park et al., 2024c), grammar/knowledge ac-
 155 quisition (Allen-Zhu & Li, 2023b;a), and interpretability
 156 methods (Menon et al., 2025; Hindupur et al., 2025). Most
 157 relevant to our work, Jain et al. (2023) used synthetic data to
 158 argue fine-tuning creates only thin wrappers over pretrained
 159 capabilities, while Nishi et al. (2024) studied formation and
 160 destruction of representational structure. However, existing
 161 synthetic frameworks typically design data generation
 162 processes without explicitly distinguishing between the un-
 163 derlying world and how data is sampled from it. Our work
 164 introduces a framework that makes this distinction explicit,

enabling systematic study of how different views of the same world shape neural representations and their downstream adaptability.

For further discussion, see App. A.

3. Experimental Framework: Decoupling World, Data and Model

Our framework uses geographic tasks where models solve geometric problems involving city coordinates. This naturally separates the underlying world (coordinates) from data generation (tasks), while providing ground-truth for measuring representation quality. Our framework provides three key properties:

1. **Learnability:** All tasks are deterministically generated from the same underlying coordinates. A model that learns the world structure can leverage it across all tasks.
2. **Latent State:** Models never see coordinates directly, only task outputs, allowing us to probe whether they internally reconstruct the world structure.
3. **Consistent Updates:** Modifying the world (e.g., adding new cities) produces self-consistent updates across all tasks, defining a clear expectation for what a model with proper world representations should internalize.

Framework. Let \mathcal{W} denote a *world*: a set of entities $\{e_1, \dots, e_N\}$ each with latent attributes $z_i \in \mathcal{Z}$. A *data generation process* is a set of tasks $\mathcal{T} = \{T_1, \dots, T_K\}$, where each task $T_k : \mathcal{Z}^{n_k} \rightarrow \mathcal{Y}_k$ maps n_k entity attributes to an output space \mathcal{Y}_k . Training data for task T_k is generated by sampling entity tuples $(e_{i_1}, \dots, e_{i_{n_k}})$ from \mathcal{W} and computing $y = T_k(z_{i_1}, \dots, z_{i_{n_k}})$.

A model M observes only entity identifiers and task outputs, never the latent attributes z_i directly. We say M has learned a *world representation* if there exists a probe P such that $P(M(e_i)) \approx z_i$ for all entities.

A *world update* $\mathcal{W} \rightarrow \mathcal{W}'$ (e.g., adding or modifying entities) induces consistent updates across all tasks by simply applying the same T_k to the new or modified entities.

Instantiation. Concretely, our world consists of 5,075 real-world cities filtered by population $> 100,000$ (Fig. 1, top). We define 7 geometric tasks that take 2 or more city coordinates as input and compute a geometric value (Fig. 1, middle).

Each task query follows a structured format where city IDs (e.g., c_1234) serve as inputs to geometric functions, all character-tokenized for autoregressive prediction. For instance, $\text{dist}(c_0865, c_4879)=769$ queries the distance between two cities, while

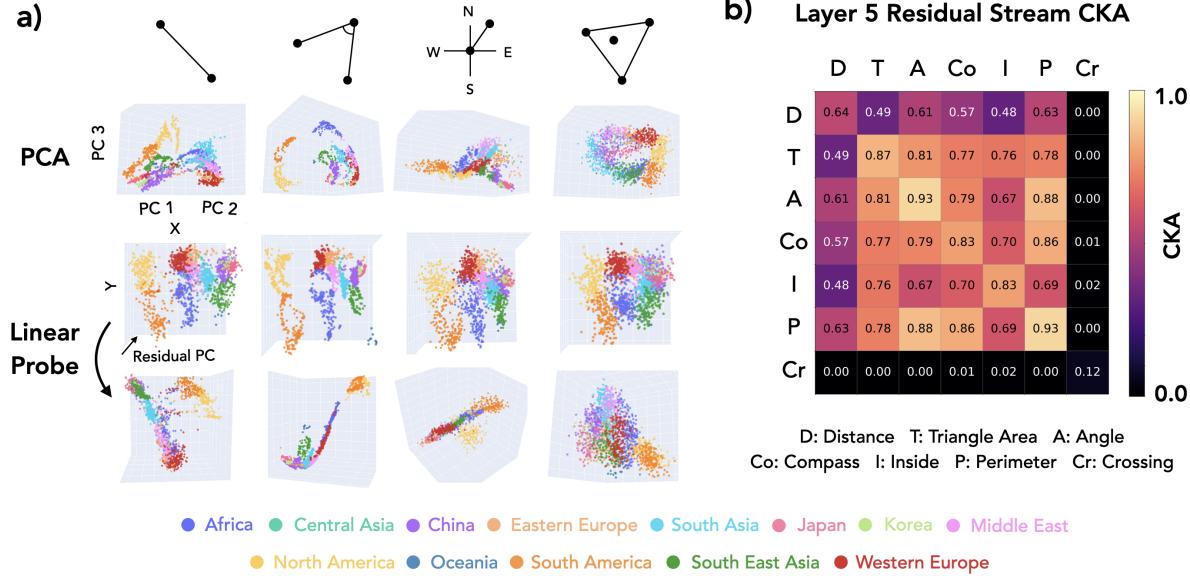


Figure 2. World representation geometry depends on the data generation process. (a) Different tasks create distinct geometries: distance (thread-like), angle (2D manifold), compass (fragmented), inside (diffuse). Row 1: PCA. Row 2: Linear probe projections. Row 3: Rotated views showing hidden structure. See App. Fig. 8 for more seeds. (b) CKA matrix at layer 5, estimated across 3 seeds. Crossing (Cr) fails to train alone. See App. Fig. 11 for SEM and layers 3, 4, 6. 3D visualizations: [link](#) (or Supp. Mat.).

cross(c_2345, c_6789; c_0123, c_4567)=TRUE
checks whether two line segments intersect.

To test adaptation, we define `Atlantis`: 100 synthetic cities placed in the Atlantic Ocean. Models never observe `Atlantis` during pretraining; we use it in Sec. 5 to test whether fine-tuning can integrate new entities into world representations in a way that generalizes across tasks.

4. World Representations Converge Under Multi-Task Learning

We now study how the task composition in the pretraining data shapes internal world representations by training Transformers on different task subsets and probing their representation geometry (see App. B.3 for training details).

Result 1: World Representations Emerge through Autoregressive Training We first demonstrate that world representations emerge through autoregressive training (Fig. 1, bottom). Training on the angle task, the model starts with random representations, goes through a loss plateau while clustering nearby cities, then forms world-aligned geometry as loss drops and task accuracy improves. The linear probe R^2 for coordinate decoding rises slightly before angle accuracy improves, reminiscent of hidden progress measures found during grokking (Nanda et al., 2023a). Notably, once representational structure forms, it remains largely fixed for the remainder of training: representations are essentially fixed in the first $\sim 15\%$ of training, remaining static while loss continues to decrease and accuracy rises (see App. 10

for visualization across tasks). This early saturation of representations echoes findings on critical learning periods in deep networks (Achille et al., 2019) and loss of plasticity in continual learning (Dohare et al., 2024). Overall, we find stable formation of internal world representations through pure autoregressive modeling. While the emergence of linearly decodable coordinates might be anticipated given the geometric nature of the task¹, it provides a useful validation of our framework and sets the stage for our main question: how do different tasks shape these representations?

Result 2: Data Generation Process Controls World Representation Geometry We train models from scratch for each of the seven tasks and visualize their representations in Fig. 2(a): PCA projections, linear probe reconstructions and rotated views.

Different tasks produce qualitatively distinct geometries: `distance` forms thread-like structures, `angle` forms 2D manifolds, `compass` forms fragmented clusters, and `inside` forms diffuse representations. These qualitative patterns are relatively consistent across random seeds (see App. D.1). Despite geometric differences, we can linearly decode (x,y) coordinates from most tasks (row 2), though some tasks (`angle`) yield cleaner reconstructions than others, a phenomenon worth further investigation. The third row shows manually rotated views revealing that representa-

¹We regard *linear* decodability of world representations as non-trivial (albeit expected). However, this is not the focus of our study.

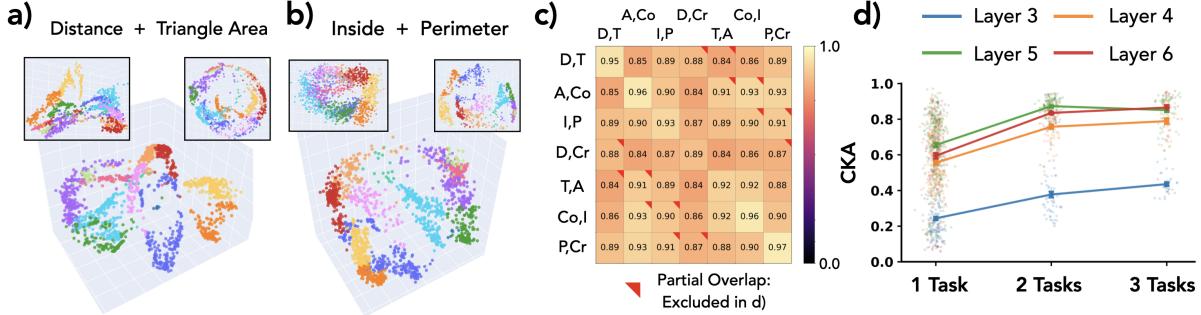


Figure 3. Multi-task pretraining drives representational convergence. (a,b) Two-task training creates more regular structures than single-task models. (c) CKA matrix (7×7) for two-task models shows higher alignment (see App. Fig. 12 for SEM). (d) Average CKA increases with task count ($1 \rightarrow 2 \rightarrow 3$), saturating at ~ 0.85 for layers 4-6 while layer 3 continues improving (see App. Fig. 13 for SEM). Crossing, which failed to learn in single-task training, is excluded; including it would only strengthen the convergence finding.

tions differ substantially in non-probe directions, a reminder that *linear probing only surfaces what we look for*.

We quantify representational similarity using CKA (Kornblith et al., 2019) (Fig. 2b). We find substantial variability even across seeds for the same task (see App. Fig. 11), but cross-task differences remain clear: distance produces particularly divergent representations, a result not obvious from PCA visualizations or from intuition about the task. Note: the crossing task fails to train in isolation², explaining its near-zero CKA; intriguingly, it succeeds in multi-task settings (Result 3).

Result 3: Multi-Task Learning Drives Representational Convergence Having established that single-task training produces variable representations, we now ask: does multi-task training reduce this variability? This question partially connects to PRH (Huh et al., 2024), which observes that neural networks trained on diverse data develop aligned representations even across different modalities and architectures. One potential mechanism they suggest is the Multitask Scaling Hypothesis:

"There are fewer representations that are competent for N tasks than there are for $M \leq N$ tasks. As we train more general models that solve more tasks at once, we should expect fewer possible solutions."

Our setup provides a potential testbed for this hypothesis, with a ground-truth world model and multiple tasks defined over it. We trained models on selected two-task combinations (3 seeds each; see App. Fig. 14 for all 21 combinations). Fig. 3(a) shows representations when trained jointly on distance and triangle area (with single-task

²This likely connects to known hard-to-learn dynamics and gradient plateaus in training transformers (Pezeshki et al., 2021; Shah et al., 2020; Hoffmann et al., 2024; Bachmann & Nagarajan, 2025; Gopalani & Hu, 2025).

models shown for comparison), while (b) shows inside and perimeter. When trained on two tasks, models develop more regular representational structures. While difficult to appreciate in static 2D projections, we encourage readers to explore our interactive 3D visualizations at [this link](#) (or Supp. Mat.).

We measure CKA between two-task trained models to quantify this alignment (Fig. 3(c)). CKA is substantially higher than for single-task models. One might expect high CKA when models share a task, but even models trained on completely disjoint task pairs show substantially higher alignment. In Fig. 3(d), we plot average CKA for models trained on 1, 2, and 3 tasks across layers 3-6, averaging only over models with completely disjoint task sets. Training on more tasks clearly leads to more aligned representations, with CKA saturating around 0.85 for 2 and 3 tasks in layers 4-6, while layer 3 continues improving. Notably, multi-task training also reduces per-seed variance of representations (App. Fig. 14b).

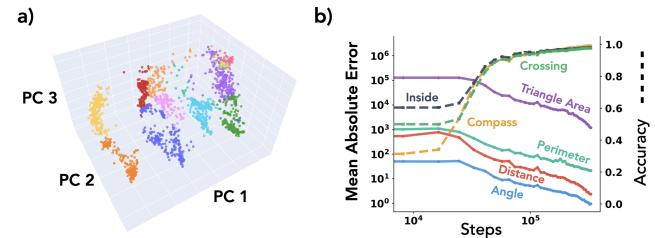


Figure 4. 7-task model. (a) PCA projection of layer 5 representations naturally reveals world map structure. (b) Training curves showing successful learning of all 7 tasks, including crossing which failed in single-task training.

Overall, we find that *multi-task learning leads to more aligned model internal representations*, providing partial evidence for the Multitask Scaling Hypothesis explanation of

PRH.³ Crucially, this alignment emerges even though single-task models achieve comparable task performance, all models reach high accuracy on their respective tasks. Since our networks are trained to representational convergence (as seen in Fig. 1), this demonstrates that the alignment is not simply a byproduct of optimization difficulty but rather that task diversity, not just data quantity or performance pressure, drives aligned representation learning.

An auxiliary finding: the `crossing` task, which was unlearnable alone, trains successfully when paired with any other task. We speculate that companion tasks provide structured coordinate representations that `crossing` can leverage, an implicit curriculum where easier tasks scaffold the learning of harder ones through shared representations.

To extend these findings, we trained a model on all 7 tasks simultaneously (Fig. 4). This model successfully learns all tasks, and its PCA projection naturally reveals the world map structure, approaching the perceived quality of linearly probed (x, y) coordinates without requiring any explicit coordinate supervision. Why multi-task training drives more linearly *surfaced* representations remains an open question worthy of future investigation. This 7-task model serves as the foundation for our fine-tuning experiments in the following section.

5. Divergent Tasks Harm Entity Integration via Fine-Tuning

In the previous section we observed how multi-task pre-training yields shared representations for different tasks. In this section, we investigate generalization properties of fine-tuning on top of such representations. However, unlike most fine-tuning studies which focus on changing model behavior in a certain way and evaluate generalization across entities, we study the inverse: fine-tuning an entity into the model and evaluate generalization across tasks. To this end, we use the 7-task model trained in the previous section (Fig. 4).

As mentioned in Sec. 3, we introduce 100 `Atlantis` cities to the world and fine-tune on data containing `Atlantis` to probe for generalization. We emphasize that the introduction of `Atlantis` cities keeps the original dataset fully consistent with the world. Moreover, task operations on `Atlantis` cities are well-defined in the same framework. If the model learned the true data generation process with properly factored representations, it should be able to integrate `Atlantis` seamlessly. If not, we suspect either the representations are fractured (Kumar et al., 2025) or gradient descent cannot trigger the right representational updates (Kumar et al., 2022).

³A full test of PRH would require showing convergence across different architectures; we test only the task-diversity mechanism here.

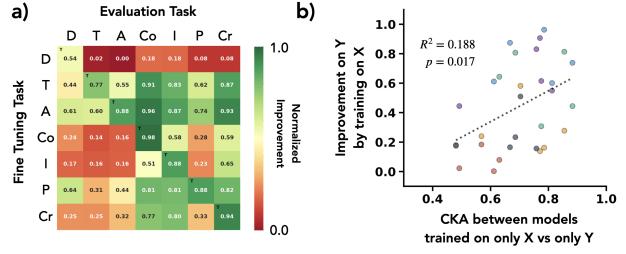


Figure 5. Fine-tuning generalization and its correlation with representational similarity. (a) Generalization matrix (averaged over 4 seeds; see App. Fig. 16 for individual seeds): each row is a model that integrated `Atlantis` via one task; columns show normalized improvement on `Atlantis` queries for each task (see App. C.1 for metric details). (b) For each task pair (X, Y), we plot the single-task CKA between X and Y against the normalized improvement on task Y after fine-tuning on task X (see App. Fig. 15 for annotated version).

Result 1: Pretraining Phase Representational Alignment Predicts Fine-Tuning Generalization Despite Joint Pretraining We first address a simple question: when fine-tuning on `Atlantis` cities for a single task (e.g., `distance`), should we expect the model to automatically generalize to using `Atlantis` for all other tasks?

To answer this, we fine-tune on 100k examples of a single task that include `Atlantis` cities, mixed with original pretraining data to avoid catastrophic forgetting and a small multi-task elicitation set (see App. B.3 for details).

The resulting generalization matrix is shown in Fig. 5(a). This matrix reveals rich phenomenology: some tasks like `distance` show no cross-task generalization (`Atlantis` remains usable only for that task), while `angle` triggers significant generalization across all tasks. Intriguingly, we observe an apparent inverse relationship: tasks that efficiently trigger cross-task generalization of new entities are often those that don't easily benefit from other tasks' fine-tuning, though this relationship is noisy.

Unexpectedly, we find that *generalization performance correlates with the CKA values from single-task pretraining* (Result 2 of Sec. 4). This is puzzling: the CKA values come from models trained from scratch on individual tasks, yet they partially predict fine-tuning behavior of a model pre-trained on all tasks jointly (Fig. 5b). If the multi-task model truly uses unified representations for cities, why would single-task representational properties matter?

For clarity, we define two terms: **Divergent tasks** are tasks which have low CKA compared to others when trained in isolation (in our case the `distance` task). **Hidden spaces** are representation spaces not surfaced by PCA or probing but used by divergent tasks.

We hypothesize:

330
 331 “Even though models develop joint world representations which converge in multi-task pretrain-
 332 ing, gradient descent on divergent tasks might fail
 333 to act on these shared representations during fine-
 334 tuning, instead utilizing hidden spaces that don’t
 335 propagate updates across tasks.”
 336

337 Our question is then two-part:

- 338 1. To what extent do divergent tasks affect fine-tuning
 339 generalization?
 340 2. Will gradient descent on divergent tasks fail to merge
 341 fine-tuning introduced concepts to the original repre-
 342 sentation manifold?

343 **Result 2: Divergent Tasks Catastrophically Harm Gener-**

344 **alization** To investigate how divergent tasks affect gener-
 345 alization, we move from single-task to multi-task fine-tuning
 346 settings. First, we introduce a simple heuristic model: fine-
 347 tuning on a concatenated dataset $\{D_1, D_2, \dots, D_n\}$ (which
 348 do not provide conflicting supervision) should combine their
 349 individual effects. Specifically, when concatenating and
 350 shuffling all fine-tuning data to avoid sequential learning
 351 effects like catastrophic forgetting (McCloskey & Cohen,
 352 1989), we expect the improvement Imp_i on task i after train-
 353 ing on tasks j and k to follow a **best-teacher model**:

354
$$\text{Imp}_i(D_j \cup D_k) = \max(\text{Imp}_i(D_j), \text{Imp}_i(D_k)) \quad (1)$$

355 To test this hypothesis, we fine-tuned the 7-task model on
 356 all $\binom{7}{2} = 21$ possible two-task combinations. Fig. 6(a,c)
 357 shows the *deviation* from our best-teacher expectation (aver-
 358 aged over 4 seeds; see App. Fig. 17 for raw improvements
 359 and App. Fig. 18 for individual seeds). Strikingly, we ob-
 360 serve “red horizontal bands”, models that not only fail to
 361 benefit from multi-task training but actually perform worse
 362 than their best single-task component. Notably, all these
 363 degraded performance bands involve the *distance* task.
 364 Fig. 6(c) quantifies this: when we split the deviation values
 365 into models with and without *distance*, we consistently
 366 observe lower-than-expected performance when the diver-
 367 gent task is included. This confirms that *divergent tasks*
 368 (*those with low single-task CKA*) *actively harm fine-tuning*
 369 *generalization rather than simply failing to contribute*. We
 370 next examine how this manifests in the learned representa-
 371 tions.

372 **Result 3: Divergent Tasks Disrupt Representational Inte-**

373 **gration of New Entities** Having shown that divergent
 374 tasks harm generalization (Question 1), we now address
 375 Question 2: does gradient descent on divergent tasks fail to
 376 merge new entities into the representation manifold?

377 We take two exemplars from the 21 fine-tuning runs:
 378 one without *distance* that generalized well (angle +

379 *compass*), and one with *distance* that was harmed
 380 (*distance + perimeter*). We first train a linear probe
 381 on a subset of all cities including *Atlantis*; these recon-
 382 structions are shown in Fig. 6(b) (top and bottom panels).
 383 In the well-integrated case, *Atlantis* cities lie within the
 384 world data manifold. In the ill-integrated case, *Atlantis*
 385 cities are off the manifold. While this difference appears
 386 subtle in 2D projections, the effect is dramatic in 3D, we
 387 strongly encourage readers to explore our [interactive visualiza-](#)
 388 [tions](#) (or Supp. Mat.). Next, we train a linear probe
 389 on 4000 non-*Atlantis* cities and apply it to *Atlantis*
 390 representations (middle panels). In the well-integrated case,
 391 *Atlantis* cities (red-orange) are relatively well recon-
 392 structed compared to ground truth (black crosses); in the
 393 ill-integrated case, reconstruction fails completely.

394 We quantify this effect in Fig. 6(d), showing histograms of
 395 absolute coordinate reconstruction error. When *Atlantis*
 396 is integrated via fine-tuning partially on divergent task data
 397 (red), reconstruction errors are nearly an order of magnitude
 398 larger than when integrated via purely non-divergent tasks
 399 (blue). For reference, non-*Atlantis* cities (yellow, still held out
 400 from probe training) show low reconstruction error as
 401 expected. One might hypothesize that *Atlantis*’s location in
 402 the middle of the ocean creates inherently difficult geometry.
 403 To test this, we pretrained a model with *Atlantis* included
 404 from the start (green line). In this case, *Atlantis* cities are
 405 reconstructed as well as any other city, confirming that the
 406 integration failure stems from divergent task fine-tuning
 407 dynamics rather than geographic peculiarity.

408 *This suggests that divergent tasks cause optimization to*
409 encode new entities in hidden spaces rather than integrating
410 them into the existing world manifold, explaining their
411 failure to support cross-task generalization.

412 We emphasize that our findings are correlational: we do not
 413 claim that interventions to increase single-task CKA would
 414 necessarily improve fine-tuning generalization. Rather, we
 415 identify representational divergence as a diagnostic marker
 416 for tasks that will harm multi-task fine-tuning performance.

417 Putting these results together: single-task representational
 418 divergence weakly predicts fine-tuning generalization even
 419 after joint pretraining, and the most divergent task (*distance*)
 420 actively harms integration of new entities. This raises a
 421 hypothesis: certain task-architecture pairings may have in-
 422 trinsic properties that induce gradient dynamics bypassing
 423 shared representations, causing updates in hidden subspaces
 424 that harm generalization, even when the network uses uni-
 425 fied representations for the forward pass.

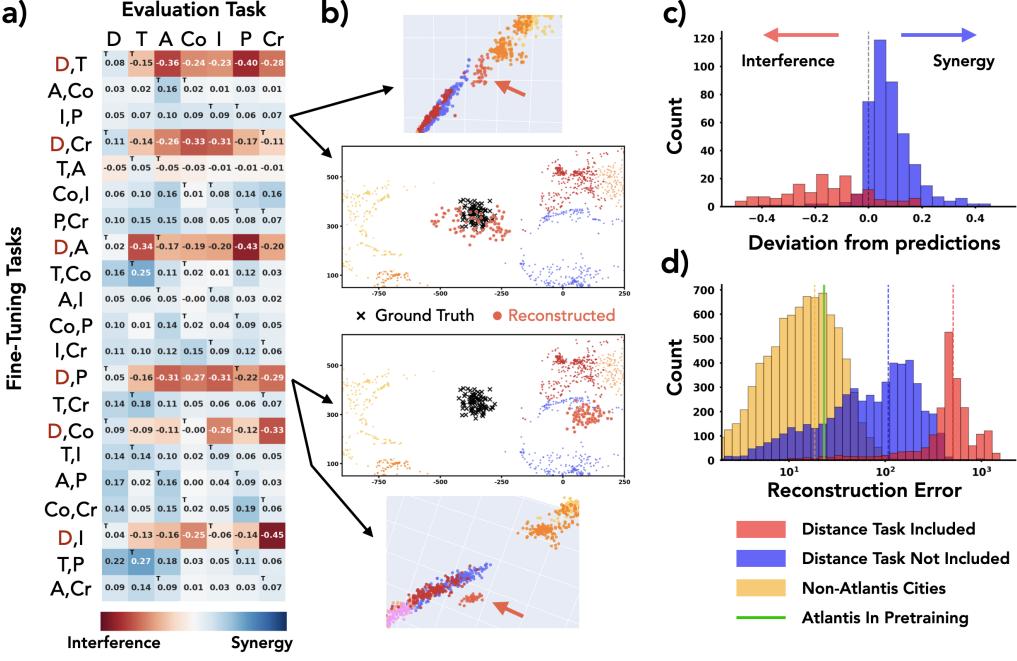


Figure 6. Divergent tasks harm multi-task fine-tuning and disrupt representational integration. (a) Deviation from best-teacher expectation for 21 two-task models (rows) across 7 evaluation tasks (columns), computed in normalized improvement space (see App. C.1); “red horizontal bands” show distance task combinations degrade performance below single-task baselines. (b) Representation visualization and linear probe reconstruction of Atlantis. (c) Histogram of deviation values: models including distance vs. not. (d) Linear probe Atlantis coordinate reconstruction error for models with *distance*, without *distance*, and baseline on pretraining cities; green vertical line indicates performance when Atlantis is part of pretraining. 3D visualizations: [link](#) (or Supp. Mat.).

6. Discussion

Continual learning and world models. For truly general intelligence, internal world models should not only represent current state but adapt consistently when the world changes. Such adaptation is non-trivial: a single change can require cascading updates across tasks. Recent language models sidestep persistent adaptation via in-context learning, forming task-specific representations on the fly (Brown et al., 2020; Park et al., 2024a; Li et al., 2025b). However, fine-tuning consistently underperforms ICL for knowledge integration (Lampinen et al., 2025; Park et al., 2025). Our study grounds these questions in a controlled setting where we can measure whether gradient descent achieves consistent integration of new entities into existing representations.

Dynamics of representations. Most recent work on neural representations examines pretrained networks or their formation during a single pretraining run. There is growing interest in how representations change during adaptation, both at inference (Park et al., 2024a; Li et al., 2025b; Shai et al., 2025; Lubana et al., 2025; Bigelow et al., 2025) and during fine-tuning (Wang et al., 2025; Minder et al., 2025; Casademunt et al., 2025). To study representational adaptation rigorously, one must define both an updatable world and how updates to it propagate into training data. Our framework provides exactly this: introducing *Atlantis*

defines how representations should update across all tasks.

Forward and backward modularity. Our results highlight a distinction that is often overlooked: *modularity in the forward pass does not imply modularity in the backward pass*. Multi-task training produces clean, structured representations that can be easily decoded into world coordinates, yet these world models can be fractured and partial when it comes to adaptation. Gradient descent may not respect the forward-pass modularity when updating weights: fine-tuning on divergent tasks routes updates through pathways that bypass the shared world manifold, encoding new entities in task-specific subspaces.

Future work. Understanding the mechanistic basis of task divergence is an important open question. If divergence is a property of task-architecture pairing rather than learned weights, it may be predictable from task structure and gradient geometry alone, enabling identification of harmful tasks before training.

Limitations. We study representation formation in a controlled synthetic setting with small-scale models; generalization to large-scale natural settings remains unclear. We identify divergence as a diagnostic marker but do not reveal underlying mechanisms. Our PRH claims are partial, as we study only a single architecture and modality.

440
441 **Impact Statement**
442

This paper presents work whose goal is to advance the field of Machine Learning. There are many potential societal consequences of our work, none which we feel must be specifically highlighted here.

443
444 **References**
445

Achille, A., Rovere, M., and Soatto, S. Critical learning periods in deep neural networks, 2019. URL <https://arxiv.org/abs/1711.08856>.

Aghajanyan, A., Gupta, A., Shrivastava, A., Chen, X., Zettlemoyer, L., and Gupta, S. Muppet: Massive multi-task representations with pre-finetuning, 2021. URL <https://arxiv.org/abs/2101.11038>.

Allen-Zhu, Z. and Li, Y. Physics of language models: Part 3.1, knowledge storage and extraction. *arXiv preprint arXiv:2309.14316*, 2023a.

Allen-Zhu, Z. and Li, Y. Physics of language models: Part 1, learning hierarchical language structures. *ArXiv e-prints, abs/2305.13673, May*, 2023b.

Anthropic AI. *Towards Monosemanticity: Decomposing Language Models With Dictionary Learning*, 2023. <https://transformer-circuits.pub/2023/monosemantic-features>.

Arditi, A., Obeso, O., Syed, A., Paleka, D., Panickssery, N., Gurnee, W., and Nanda, N. Refusal in language models is mediated by a single direction. *arXiv preprint arXiv:2406.11717*, 2024.

Bachmann, G. and Nagarajan, V. The pitfalls of next-token prediction, 2025. URL <https://arxiv.org/abs/2403.06963>.

Bengio, Y., Courville, A., and Vincent, P. Representation learning: A review and new perspectives, 2014. URL <https://arxiv.org/abs/1206.5538>.

Berglund, L., Tong, M., Kaufmann, M., Balesni, M., Stickland, A. C., Korbak, T., and Evans, O. The reversal curse: Llms trained on "a is b" fail to learn "b is a", 2024. URL <https://arxiv.org/abs/2309.12288>.

Betley, J., Tan, D., Warncke, N., Sztyber-Betley, A., Bao, X., Soto, M., Labenz, N., and Evans, O. Emergent misalignment: Narrow finetuning can produce broadly misaligned llms, 2025. URL <https://arxiv.org/abs/2502.17424>.

Bigelow, E., Wurgaft, D., Wang, Y., Goodman, N., Ullman, T., Tanaka, H., and Lubana, E. S. Belief dynamics reveal

the dual nature of in-context learning and activation steering, 2025. URL <https://arxiv.org/abs/2511.00617>.

Bronstein, M. M., Bruna, J., Cohen, T., and Veličković, P. Geometric deep learning: Grids, groups, graphs, geodesics, and gauges, 2021. URL <https://arxiv.org/abs/2104.13478>.

Brown, T., Mann, B., Ryder, N., Subbiah, M., Kaplan, J. D., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33: 1877–1901, 2020.

Caruana, R. Multitask learning. *Machine learning*, 28(1): 41–75, 1997.

Casademunt, H., Juang, C., Karvonen, A., Marks, S., Rajamanoharan, S., and Nanda, N. Steering out-of-distribution generalization with concept ablation finetuning, 2025. URL <https://arxiv.org/abs/2507.16795>.

Chan, S. C. Y., Santoro, A., Lampinen, A. K., Wang, J. X., Singh, A., Richemond, P. H., McClelland, J., and Hill, F. Data distributional properties drive emergent in-context learning in transformers, 2022. URL <https://arxiv.org/abs/2205.05055>.

Cohen, T. S. and Welling, M. Group equivariant convolutional networks, 2016. URL <https://arxiv.org/abs/1602.07576>.

Demircan, C., Saanum, T., Jagadish, A. K., Binz, M., and Schulz, E. Sparse autoencoders reveal temporal difference learning in large language models, 2024. URL <https://arxiv.org/abs/2410.01280>.

Devlin, J., Chang, M.-W., Lee, K., and Toutanova, K. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.

Dohare, S., Hernandez-Garcia, J. F., Rahman, P., Mahmood, A. R., and Sutton, R. S. Maintaining plasticity in deep continual learning, 2024. URL <https://arxiv.org/abs/2306.13812>.

Fu, S., Bonnen, T., Guillory, D., and Darrell, T. Hidden in plain sight: Vlms overlook their visual representations, 2025. URL <https://arxiv.org/abs/2506.08008>.

Fukushima, K. Neocognitron: A self-organizing neural network model for a mechanism of pattern recognition unaffected by shift in position. *Biological cybernetics*, 36(4):193–202, 1980.

- 495 Ge, X., Shu, W., Wu, J., Zhou, Y., He, Z., and Qiu, X. Evolution
 496 of concepts in language model pre-training, 2025.
 497 URL <https://arxiv.org/abs/2509.17196>.
- 498 Gopalani, P. and Hu, W. What happens during the loss
 499 plateau? understanding abrupt learning in transfor-
 500 mers, 2025. URL <https://arxiv.org/abs/2506.13688>.
- 501 502
- 503 Gurnee, W. and Tegmark, M. Language models represent
 504 space and time. *arXiv preprint arXiv:2310.02207*, 2023.
- 505 506 He, K., Zhang, X., Ren, S., and Sun, J. Deep residual
 507 learning for image recognition, 2015.
- 508 Higgins, I., Matthey, L., Pal, A., Burgess, C., Glorot, X.,
 509 Botvinick, M., Mohamed, S., and Lerchner, A. betavae:
 510 Learning basic visual concepts with a constrained
 511 variational framework. In *Proc. Int. Conf. on Learning
 512 Representations (ICLR)*, 2017.
- 513 514 Hindupur, S. S. R., Lubana, E. S., Fel, T., and Ba, D.
 515 Projecting assumptions: The duality between sparse au-
 516 toencoders and concept geometry, 2025. URL <https://arxiv.org/abs/2503.01822>.
- 517 518
- 519 Hoffmann, D. T., Schrodi, S., Bratulić, J., Behrmann, N.,
 520 Fischer, V., and Brox, T. Eureka-moments in transform-
 521 ers: Multi-step tasks reveal softmax induced optimization
 522 problems, 2024. URL <https://arxiv.org/abs/2310.12956>.
- 523 524
- 525 Hu, E. J., Shen, Y., Wallis, P., Allen-Zhu, Z., Li, Y., Wang,
 526 S., Wang, L., and Chen, W. Lora: Low-rank adaptation of
 527 large language models, 2021. URL <https://arxiv.org/abs/2106.09685>.
- 528 529
- 530 Hubel, D. H. and Wiesel, T. N. Receptive fields, binocular
 531 interaction and functional architecture in the cat's visual
 532 cortex. *The Journal of physiology*, 160(1):106, 1962.
- 533 534
- 535 Huh, M., Cheung, B., Wang, T., and Isola, P. The pla-
 536 tonic representation hypothesis, 2024. URL <https://arxiv.org/abs/2405.07987>.
- 537 538
- 539 Ilharco, G., Ribeiro, M. T., Wortsman, M., Gururangan,
 540 S., Schmidt, L., Hajishirzi, H., and Farhadi, A. Editing
 541 models with task arithmetic, 2023. URL <https://arxiv.org/abs/2212.04089>.
- 542 543
- 544 Jain, S., Kirk, R., Lubana, E. S., Dick, R. P., Tanaka,
 545 H., Grefenstette, E., Rocktäschel, T., and Krueger,
 546 D. S. Mechanistically analyzing the effects of fine-
 547 tuning on procedurally defined tasks. *arXiv preprint
 548 arXiv:2311.12786*, 2023.
- 549 550
- 551 Kim, J., Kwon, S., Choi, J. Y., Park, J., Cho, J., Lee,
 552 J. D., and Ryu, E. K. Task diversity shortens the icl
 553 554 plateau, 2025. URL <https://arxiv.org/abs/2410.05448>.
- 555 556
- 557 Kornblith, S., Norouzi, M., Lee, H., and Hinton, G. Similar-
 558 ity of Neural Network Representations Revisited. In *Proc.
 559 of the 36th Proc. Int. Conf. on Machine Learning (ICML)*,
 560 Proc. of Machine Learning Research. PMLR, 09–15 Jun
 561 2019.
- 562 563
- 564 Krizhevsky, A., Sutskever, I., and Hinton, G. E. Imagenet
 565 classification with deep convolutional neural networks.
 566 *Advances in neural information processing systems*, 25,
 567 2012.
- 568 569
- 570 Kumar, A., Raghunathan, A., Jones, R., Ma, T., and Liang,
 571 P. Fine-tuning can distort pretrained features and un-
 572 derperform out-of-distribution, 2022. URL <https://arxiv.org/abs/2202.10054>.
- 573 574
- 575 Kumar, A., Clune, J., Lehman, J., and Stanley, K. O. Ques-
 576 tioning representational optimism in deep learning: The
 577 fractured entangled representation hypothesis, 2025. URL
 578 <https://arxiv.org/abs/2505.11581>.
- 579 580
- 581 Lampinen, A. K., Chaudhry, A., Chan, S. C. Y., Wild, C.,
 582 Wan, D., Ku, A., Bornschein, J., Pascanu, R., Shanahan,
 583 M., and McClelland, J. L. On the generalization
 584 of language models from in-context learning and
 585 finetuning: a controlled study, 2025. URL <https://arxiv.org/abs/2505.00661>.
- 586 587
- 588 Lee, A., Bai, X., Pres, I., Wattenberg, M., Kummerfeld,
 589 J. K., and Mihalcea, R. A mechanistic understanding of
 590 alignment algorithms: A case study on dpo and toxic-
 591 ity. In *Forty-first International Conference on Machine
 592 Learning*, 2024. URL <https://arxiv.org/abs/2401.01967>.
- 593 594
- 595 Lee, A., Sun, L., Wendler, C., Viégas, F., and Wattenberg,
 596 M. The geometry of self-verification in a task-specific
 597 reasoning model, 2025. URL <https://arxiv.org/abs/2504.14379>.
- 598 599
- 600 Lester, B., Al-Rfou, R., and Constant, N. The power of
 601 scale for parameter-efficient prompt tuning, 2021. URL
 602 <https://arxiv.org/abs/2104.08691>.
- 603 604
- 605 Li, K., Hopkins, A. K., Bau, D., Viégas, F., Pfister, H.,
 606 and Wattenberg, M. Emergent world representations:
 607 Exploring a sequence model trained on a synthetic task.
 608 In *The Eleventh International Conference on Learning
 609 Representations*, 2022.
- 610 611
- 612 Li, M. Z., Agrawal, K. K., Ghosh, A., Teru, K. K., La-
 613 joie, G., and Richards, B. A. Tracing the representation
 614 geometry of language models from pretraining to post-
 615 training. In *High-dimensional Learning Dynamics 2025,
 616 2025a*. URL <https://openreview.net/forum?id=9nKmDLXg9v>.

- 550 Li, Y., Campbell, D., Chan, S. C. Y., and Lampinen, A. K.
551 Just-in-time and distributed task representations in lan-
552 guage models, 2025b. URL <https://arxiv.org/abs/2509.04466>.
- 553
- 554 Lubana, E. S., Rager, C., Hindupur, S. S. R., Costa, V.,
555 Tuckute, G., Patel, O., Murthy, S. K., Fel, T., Wurgaft, D.,
556 Bigelow, E. J., Lin, J., Ba, D., Wattenberg, M., Viegas,
557 F., Weber, M., and Mueller, A. Priors in time: Missing
558 inductive biases for language model interpretability, 2025.
559 URL <https://arxiv.org/abs/2511.01836>.
- 560
- 561 Malladi, S., Gao, T., Nichani, E., Damian, A., Lee, J. D.,
562 Chen, D., and Arora, S. Fine-tuning language models
563 with just forward passes, 2024. URL <https://arxiv.org/abs/2305.17333>.
- 564
- 565 Marks, S. and Tegmark, M. The geometry of truth: Emer-
566 gent linear structure in large language model represen-
567 tations of true/false datasets, 2024. URL <https://arxiv.org/abs/2310.06824>.
- 568
- 569 McCloskey, M. and Cohen, N. J. Catastrophic interfer-
570 ence in connectionist networks: The sequential learning
571 problem. In *Psychology of learning and motivation*, vol-
572 ume 24, pp. 109–165. Elsevier, 1989.
- 573
- 574 Menon, A., Shrivastava, M., Krueger, D., and Lubana, E. S.
575 Analyzing (in)abilities of saes via formal languages, 2025.
576 URL <https://arxiv.org/abs/2410.11767>.
- 577
- 578 Michaud, E. J., Liu, Z., Girit, U., and Tegmark, M. The
579 quantization model of neural scaling. *arXiv preprint*
580 *arXiv:2303.13506*, 2023.
- 581
- 582 Minder, J., Dumas, C., Juang, C., Chugtai, B., and Nanda,
583 N. Overcoming sparsity artifacts in crosscoders to inter-
584 pret chat-tuning, 2025. URL <https://arxiv.org/abs/2504.02922>.
- 585
- 586 Mircea, A., Chakraborty, S., Chitsazan, N., Naphade, M.,
587 Sahu, S., Rish, I., and Lobacheva, E. Training dynamics
588 underlying language model scaling laws: Loss de-
589 celeration and zero-sum learning, 2025. URL <https://arxiv.org/abs/2506.05447>.
- 590
- 591 Nanda, N., Chan, L., Lieberum, T., Smith, J., and Stein-
592 hardt, J. Progress measures for grokking via mechanistic
593 interpretability, 2023a. URL <https://arxiv.org/abs/2301.05217>.
- 594
- 595 Nanda, N., Lee, A., and Wattenberg, M. Emergent lin-
596 ear representations in world models of self-supervised
597 sequence models. In *Proceedings of the 6th Black-
598 boxNLP Workshop: Analyzing and Interpreting Neural
599 Networks for NLP*, pp. 16–30, 2023b. URL <https://arxiv.org/abs/2309.00941>.
- 600
- 601 Nishi, K., Okawa, M., Ramesh, R., Khona, M., Lubana,
602 E. S., and Tanaka, H. Representation shattering in trans-
603 formers: A synthetic study with knowledge editing. *arXiv*
604 *preprint arXiv:2410.17194*, 2024.
- 605
- 606 Okawa, M., Lubana, E. S., Dick, R. P., and Tanaka, H. Com-
607 positional abilities emerge multiplicatively: Exploring
608 diffusion models on a synthetic task, 2024.
- 609
- 610 Olah, C., Mordvintsev, A., and Schubert, L. Feature vi-
611 sualization. *Distill*, 2017. doi: 10.23915/distill.00007.
612 <https://distill.pub/2017/feature-visualization>.
- 613
- 614 OpenDataSoft / GeoNames. Geonames – all cities
615 with a population ≥ 1000 . <https://public.opendatasoft.com/explore/dataset/geonames-all-cities-with-a-population-1000>,
616 2025. Accessed: 2025.
- 617
- 618 Park, C. F., Lee, A., Lubana, E. S., Yang, Y., Okawa, M.,
619 Nishi, K., Wattenberg, M., and Tanaka, H. Iclr: In-context
620 learning of representations, 2024a. URL <https://arxiv.org/abs/2501.00070>.
- 621
- 622 Park, C. F., Lubana, E. S., Pres, I., and Tanaka, H. Compe-
623 tition dynamics shape algorithmic phases of in-context
624 learning. *arXiv preprint arXiv:2412.01003*, 2024b.
- 625
- 626 Park, C. F., Okawa, M., Lee, A., Lubana, E. S., and Tanaka,
627 H. Emergence of hidden capabilities: Exploring learning
628 dynamics in concept space, 2024c. URL <https://arxiv.org/abs/2406.19370>.
- 629
- 630 Park, C. F., Zhang, Z., and Tanaka, H. New News: System-2
631 fine-tuning for robust integration of new knowledge, 2025.
632 URL <https://arxiv.org/abs/2505.01812>.
- 633
- 634 Pearce, M., Simon, E., Byun, M., and Balsam, D. Finding
635 the tree of life in evo 2. *Goodfire Research*, August 2025.
636 Correspondence to michael@goodfire.ai.
- 637
- 638 Pezeshki, M., Kaba, O., Bengio, Y., Courville, A. C., Pre-
639 cup, D., and Lajoie, G. Gradient starvation: A learning
640 proclivity in neural networks. *Adv. in Neural Information
641 Processing Systems (NeurIPS)*, 2021.
- 642
- 643 Qin, T., Park, C. F., Kwun, M., Walsman, A., Malach, E.,
644 Anand, N., Tanaka, H., and Alvarez-Melis, D. Decom-
645 posing elements of problem solving: What “math” does
646 rl teach?, 2025. URL <https://arxiv.org/abs/2505.22756>.
- 647
- 648 Radford, A., Narasimhan, K., Salimans, T., Sutskever, I.,
649 et al. Improving language understanding by generative
650 pre-training, 2018.
- 651
- 652 Raventós, A., Paul, M., Chen, F., and Ganguli, S. Pretrain-
653 ing task diversity and the emergence of non-bayesian

- 605 in-context learning for regression, 2023. URL <https://arxiv.org/abs/2306.15063>.
- 606
- 607 Reddy, G. The mechanistic basis of data dependence and
- 608 abrupt learning in an in-context classification task, 2023.
- 609 URL <https://arxiv.org/abs/2312.03002>.
- 610
- 611 Rosenblatt, F. The perceptron: a probabilistic model for
- 612 information storage and organization in the brain. *Psychological review*, 65(6):386, 1958.
- 613
- 614 Rumelhart, D. E., Hinton, G. E., and Williams, R. J. Learning
- 615 representations by back-propagating errors. *nature*, 323(6088):533–536, 1986.
- 616
- 617 Shah, H., Tamuly, K., Raghunathan, A., Jain, P., and Ne-
- 618 trapalli, P. The pitfalls of simplicity bias in neural
- 619 networks, 2020. URL <https://arxiv.org/abs/2006.07710>.
- 620
- 621 Shai, A. S., Marzen, S. E., Teixeira, L., Oldenziel, A. G.,
- 622 and Riechers, P. M. Transformers represent belief state
- 623 geometry in their residual stream, 2025. URL <https://arxiv.org/abs/2405.15943>.
- 624
- 625 Singh, A. K., Moskovitz, T., Hill, F., Chan, S. C. Y., and
- 626 Saxe, A. M. What needs to go right for an induction head?
- 627 a mechanistic study of in-context learning circuits and
- 628 their formation, 2024. URL <https://arxiv.org/abs/2404.07129>.
- 629
- 630 Templeton, A., Conerly, T., Marcus, J., Lindsey, J., Bricken,
- 631 T., Chen, B., Pearce, A., Citro, C., Ameisen, E., Jones,
- 632 A., Cunningham, H., Turner, N. L., McDougall, C.,
- 633 MacDiarmid, M., Freeman, C. D., Sumers, T. R.,
- 634 Rees, E., Batson, J., Jermyn, A., Carter, S., Olah,
- 635 C., and Henighan, T. Scaling monosemanticity: Ex-
- 636 tracting interpretable features from claude 3 sonnet.
- 637 *Transformer Circuits Thread*, 2024. URL <https://transformer-circuits.pub/2024/scaling-monosemanticity/index.html>.
- 638
- 639 Treutlein, J., Choi, D., Betley, J., Marks, S., Anil, C., Grosse,
- 640 R., and Evans, O. Connecting the dots: Llms can in-
- 641 fer and verbalize latent structure from disparate training
- 642 data, 2024. URL <https://arxiv.org/abs/2406.14546>.
- 643
- 644 Vafa, K., Chang, P. G., Rambachan, A., and Mullainathan,
- 645 S. What has a foundation model found? using inductive
- 646 bias to probe for world models, 2025. URL <https://arxiv.org/abs/2507.06952>.
- 647
- 648 Wang, A., Engels, J., Clive-Griffin, O., Rajamanoharan, S.,
- 649 and Nanda, N. Simple mechanistic explanations for out-
- 650 of-context reasoning, 2025. URL <https://arxiv.org/abs/2507.08218>.
- 651
- 652 Ward, J., Lin, C., Venhoff, C., and Nanda, N. Reasoning-
- 653 finetuning repurposes latent representations in base mod-
- 654 els, 2025. URL <https://arxiv.org/abs/2507.12638>.
- 655
- Weiler, M. and Cesa, G. General $e(2)$ -equivariant steer-
- 656 able cnns, 2021. URL <https://arxiv.org/abs/1911.08251>.
- 657
- Wu, Z., Arora, A., Wang, Z., Geiger, A., Jurafsky, D.,
- 658 Manning, C. D., and Potts, C. Reft: Representation
- 659 finetuning for language models, 2024. URL <https://arxiv.org/abs/2404.03592>.
- 660
- Wurgaft, D., Lubana, E. S., Park, C. F., Tanaka, H., Reddy,
- 661 G., and Goodman, N. D. In-context learning strategies
- 662 emerge rationally, 2025. URL <https://arxiv.org/abs/2506.17859>.
- 663
- Xie, S. M., Raghunathan, A., Liang, P., and Ma, T. An
- 664 explanation of in-context learning as implicit bayesian
- 665 inference. *arXiv preprint arXiv:2111.02080*, 2021.
- 666
- Yue, Y., Chen, Z., Lu, R., Zhao, A., Wang, Z., Yue, Y.,
- 667 Song, S., and Huang, G. Does reinforcement learning
- 668 really incentivize reasoning capacity in llms beyond
- 669 the base model?, 2025. URL <https://arxiv.org/abs/2504.13837>.
- 670
- Zhang, S., Patel, A., Rizvi, S. A., Liu, N., He, S., Karbasi,
- 671 A., Zappala, E., and van Dijk, D. Intelligence at the edge
- 672 of chaos, 2025. URL <https://arxiv.org/abs/2410.02536>.
- 673
- Zhao, R., Meterez, A., Kakade, S., Pehlevan, C., Jelassi, S.,
- 674 and Malach, E. Echo chamber: RI post-training amplifies
- 675 behaviors learned in pretraining, 2025. URL <https://arxiv.org/abs/2504.07912>.
- 676
- Zweiger, A., Pari, J., Guo, H., Akyürek, E., Kim, Y., and
- 677 Agrawal, P. Self-adapting language models, 2025. URL
- 678 <https://arxiv.org/abs/2506.10943>.
- 679

660

661

662

663

664

665

666

667

668

669

670

671

672

673

674

675

676

677

Appendix

A. Extended Related Work

Internal Representations. Understanding internal representations has roots in neuroscience (Hubel & Wiesel, 1962), informing early neural network development (Fukushima, 1980; Bengio et al., 2014; Rosenblatt, 1958; Rumelhart et al., 1986). Recent work has revealed that language models develop structured “world models” encoding geographic, temporal and relational information (Li et al., 2022; Gurnee & Tegmark, 2023; Nanda et al., 2023b; Marks & Tegmark, 2024), with similar representations emerging during in-context learning (Vafa et al., 2025). Mechanistic interpretability and sparse autoencoders have enabled decomposition of neural activations into interpretable features (Anthropic AI, 2023; Templeton et al., 2024). Researchers have also uncovered that models represent meaningful properties of data—concepts (Pearce et al., 2025; Higgins et al., 2017), features (Olah et al., 2017), and abstractions (Lee et al., 2025; Arditì et al., 2024)—in interpretable ways. Furthermore, PRH posits that diverse models converge toward similar representational structures (Huh et al., 2024). However, recent work questions this representational optimism, suggesting that deep network representations may be more brittle than previously assumed (Kumar et al., 2025). Only recent work has begun examining how representations emerge during pretraining in real LLMs (Li et al., 2025a; Ge et al., 2025) or how they change during fine-tuning (Lee et al., 2024). Our work takes a complementary perspective, studying the factors that control the formation of these representations and how networks integrate new entities into their representation space via fine-tuning.

Fine-tuning. The pretraining-finetuning paradigm has become central to modern deep learning, with seminal works establishing its effectiveness in computer vision (Krizhevsky et al., 2012; He et al., 2015) and natural language processing (Devlin et al., 2018; Radford et al., 2018). Despite widespread success, fine-tuning exhibits poorly understood behaviors such as the reversal curse (Berglund et al., 2024; Lampinen et al., 2025), out-of-context reasoning limitations (Treutlein et al., 2024), and off-target effects (Betley et al., 2025). On this background, careful studies of fine-tuning and other low-compute adaptation methods have raised pessimism about whether models can learn fundamentally new abilities, suggesting they may merely form “thin wrappers” around pretrained representations (Jain et al., 2023; Ward et al., 2025; Yue et al., 2025; Qin et al., 2025; Zhao et al., 2025; Zweiger et al., 2025). Fine-tuning has also been studied across diverse directions: parameter efficiency (Hu et al., 2021; Lester et al., 2021), zeroth-order optimization (Malladi et al., 2024), weight composition (Ilharco et al., 2023), and representation adaptation (Wu et al., 2024). Work on feature distortion (Kumar et al., 2022) is perhaps most related to ours, though representational changes are assumed rather than directly measured. Our work examines this question in a controlled setup where ground-truth world structure enables precise measurement of representation adaptation.

Multi-task Learning. Multi-task learning has long been studied as a way to improve generalization through shared representations (Caruana, 1997); in some sense, modern foundation models represent an extreme form of multi-task training. Large-scale multi-task pretraining typically assumes rich representations emerge from data diversity (Aghajanyan et al., 2021), but the precise mechanisms remain underexplored. Recent work has begun studying task diversity in controlled settings (Michaud et al., 2023; Zhang et al., 2025), though most studies still focus on aggregate behaviors such as capabilities and scaling laws rather than characterizing tasks or the knowledge they operate on. Our framework explicitly defines tasks as geometric functions over a shared world, enabling direct investigation of how task structure shapes representations.

Synthetic Data. The cost and complexity of foundation models has motivated synthetic approaches for controlled study of in-context learning (Xie et al., 2021; Chan et al., 2022; Reddy, 2023; Roventós et al., 2023; Park et al., 2024b; Wurgaft et al., 2025), compositional generalization (Okawa et al., 2024; Park et al., 2024c), grammar/knowledge acquisition (Allen-Zhu & Li, 2023b;a), and interpretability methods (Menon et al., 2025; Hindupur et al., 2025). Most relevant to our work, Jain et al. (2023) used synthetic data to argue fine-tuning creates only thin wrappers over pretrained capabilities, while Nishi et al. (2024) studied formation and destruction of representational structure. However, existing synthetic frameworks typically design data generation processes without explicitly distinguishing between the underlying world and how data is sampled from it. Our work introduces a framework that makes this distinction explicit, enabling systematic study of how different views of the same world shape neural representations and their downstream adaptability.

Dynamics of Representations. Recent work has begun studying how representations evolve during in-context learning (Shai et al., 2025; Demircan et al., 2024) or fine-tuning (Casademunt et al., 2025; Minder et al., 2025). Relatedly, Lubana et al. (2025) show that representations exhibit rich temporal dynamics that standard interpretability methods (e.g., SAEs) fail to capture due to stationarity assumptions. Fu et al. (2025) show that VLMs trained by merging LLMs and vision encoders often fail to utilize representations surfaced by the vision encoder, i.e. the representations exist but remain unused.

Geometric Deep Learning. Geometric deep learning studies how data geometry interacts with model architectures,

713

714

developing equivariant networks that respect symmetries (Bronstein et al., 2021; Cohen & Welling, 2016; Weiler & Cesa, 2021). While our world is defined on a 2D plane, one might ask: why not a sphere, torus, or other manifold? This is an interesting direction, but not our focus. We study how neural networks adapt internal representations to tasks in an arbitrarily chosen geometry. Moreover, a change in world geometry can be absorbed into the task definition (e.g., geodesic vs. Euclidean distance), so the key question remains how representations form given the task, not the underlying manifold. Planar coordinates also allow clean linear probing of world representations. Our models are standard transformers without geometric priors; we study what representations emerge purely from training on task data, treating geometry as emergent rather than imposed.

Loss Plateaus. Our crossing task fails to learn in single-task training despite escaping an initial plateau (likely output format learning), suggesting it remains stuck in a deeper plateau. Such plateaus are notoriously difficult for transformers. Recent work has studied this phenomenon mechanistically in transformers (Hoffmann et al., 2024; Gopalani & Hu, 2025; Singh et al., 2024), while others relate it to more general optimization challenges in deep learning such as simplicity bias and gradient starvation (Shah et al., 2020; Pezeshki et al., 2021; Bachmann & Nagarajan, 2025). Most related to our findings, Kim et al. (2025) show that multi-task training shortens loss plateaus, similar to why our crossing task trains successfully when joined with any other task.

B. Experimental Details

This section provides detailed information about the world, data generation process, model architecture and training procedures used in our experiments.

B.1. World

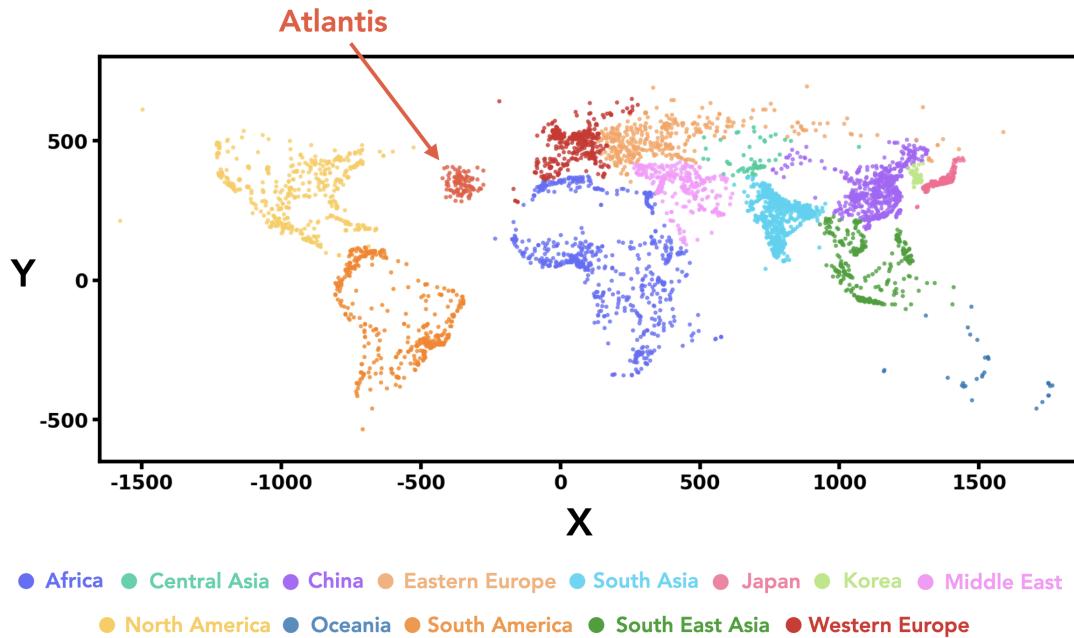


Figure 7. **Geographic distribution of cities used in our experiments.** 5,075 real-world cities plus 100 synthetic Atlantis cities (5,175 total). Cities span all continents and provide a fixed, measurable world structure. Coordinates use an equirectangular projection: $x = 10 \times \text{longitude}$, $y = 10 \times \text{latitude}$ (in degrees). The Atlantis region (Atlantic Ocean) is used for out-of-distribution testing.

Our experiments use a geographic world consisting of 5,075 cities extracted from the GeoNames database (OpenDataSoft / GeoNames, 2025) with population greater than 100,000. Cities are distributed across all continents. This choice provides natural variation in density (e.g., dense regions like India versus sparse Oceania) that creates interesting computational challenges.

While we use real city coordinates, this work studies abstract geometric reasoning rather than actual geography, we project coordinates to Euclidean space using an equirectangular projection (as described above) and treat all tasks as pure geometry

770 problems.

771 Additionally, we introduce 100 synthetic *Atlantis* cities positioned in the Atlantic Ocean, centered at (longitude -35° ,
 772 latitude 35°) and following a Gaussian distribution with standard deviation of 3° . These synthetic cities enable controlled
 773 out-of-distribution experiments, as models never observe *Atlantis* during pretraining but must generalize to it during
 774 evaluation. City IDs are randomly assigned from the range [0, 9999], creating a sparse identifier space that models must
 775 learn to map to coordinates.
 776

777 B.2. Data Generation Process

779 **Tasks** We implement 7 geometric tasks that operate on city coordinates. All tasks use a consistent format:
 780 `task(arguments)=answer`, where city IDs are prefixed with `c_`. Numerical outputs (distance, area, angle, perimeter)
 781 are rounded to integers. Table 1 summarizes the tasks.
 782

Task	Input	Output Type	Unit/Values	Example
distance	2 cities	Numerical	Scaled coords	<code>dist(c_865, c_4879)=769</code>
triarea	3 cities	Numerical	Scaled coords ²	<code>triarea(c_1234, c_5678, c_9012)=45823</code>
angle	3 cities	Numerical	Degrees (0–180)	<code>angle(c_2345, c_6789, c_123)=97</code>
compass	2 cities	Categorical	8 directions	<code>compass(c_1234, c_5678)=NE</code>
inside	1 + n cities	Categorical	TRUE/FALSE	<code>inside(c_9012; c_3456, ...)=FALSE</code>
perimeter	n cities	Numerical	Scaled coords	<code>perimeter(c_4567, c_8901, ...)=2856</code>
crossing	4 cities	Categorical	TRUE/FALSE	<code>cross(c_2345, c_6789; c_123, c_4567)=TRUE</code>

791 *Table 1.* Summary of 7 geometric tasks. Numerical outputs are integers; “scaled coords” refers to the $\times 10$ coordinate system (Sec. B.1).
 792 Categorical tasks have discrete outputs: `compass` uses 8 cardinal directions (N, NE, E, SE, S, SW, W, NW), while `inside` and
 793 `crossing` are binary. The `inside` task tests if the first city lies within the convex hull of the remaining cities; `crossing` tests if line
 794 segment (c_1, c_2) intersects segment (c_3, c_4) .

795
 796
 797 **Dataset Sizes** Each pretraining set consists of 1M rows of data per task. For fine-tuning, the dataset consists of: (1) 100k
 798 rows of the target task containing at least one *Atlantis* city, (2) 20k rows randomly sampled from the original pretraining
 799 data to prevent catastrophic forgetting, and (3) 256 rows per task (without *Atlantis*) to elicit multi-task performance.
 800

801 B.3. Model and Training

803 **Tokenization** We use character-level tokenization with 98 ASCII tokens (excluding space, which serves as the delimiter),
 804 plus special tokens for beginning-of-sequence (BOS), end-of-sequence (EOS) and padding (PAD).
 805

806 **Architecture** We use the Qwen2 decoder-only transformer architecture with hidden size 128, 4 attention heads and 6
 807 layers.
 808

809 **Pretraining** We train models autoregressively on the full sequence (no prompt masking). All pretraining runs see 42M rows
 810 regardless of dataset size (e.g., 42 epochs for 1M rows, 6 epochs for 7M rows). Table 2 summarizes the hyperparameters.
 811

Hyperparameter	Value
Optimizer	AdamW
Learning rate	3×10^{-4}
Weight decay	0.01
Scheduler	Linear with warmup
Warmup steps	50
Batch size	128
Max sequence length	256
Total training rows	42M
Initialization scale	0.1 (std)

822 *Table 2. Pretraining hyperparameters.*
 823
 824

825 **Fine-Tuning** Fine-tuning starts from the final pretrained checkpoint. We use a reduced learning rate of 1×10^{-5} ($30\times$
 826 smaller than pretraining) to avoid catastrophic forgetting. The fine-tuning dataset consists of 100k rows per task containing
 827 at least one `Atlantis` city. We train for 30 epochs with batch size 128.
 828

829 C. Analysis Methods

830 C.1. Evaluation

831 **Generation Protocol** For evaluation, we use teacher forcing up to the “=” sign (the prompt), then generate autoregressively
 832 at temperature zero until reaching the EOS token or a maximum of 128 tokens (sufficient for all tasks).
 833

834 **Task-Specific Metrics** Categorical tasks (`compass`, `inside`, `crossing`) are evaluated using accuracy. Numerical
 835 tasks are evaluated using absolute error.
 836

837 **Normalized Improvement** To compare generalization across tasks with different metrics and scales, we define a
 838 normalized improvement score that maps performance to $[0, 1]$, where 0 indicates no improvement over the `Atlantis`
 839 baseline (before fine-tuning) and 1 indicates matching the pretrained model’s performance on standard cities.
 840

841 For **error-based tasks** (`distance`, `triarea`, `angle`, `perimeter`), where lower is better:

$$842 \quad NI = \frac{\log(\text{baseline}_{\text{atlantis}}/\text{error})}{\log(\text{baseline}_{\text{atlantis}}/\text{baseline}_{\text{standard}})} \quad (2)$$

843 The logarithmic scaling ensures multiplicative improvements are treated equally (e.g., reducing error from 1000 to 100 is
 844 weighted the same as 100 to 10).
 845

846 For **accuracy-based tasks** (`compass`, `inside`, `crossing`), where higher is better:

$$847 \quad NI = \frac{\text{accuracy} - \text{baseline}_{\text{atlantis}}}{\text{baseline}_{\text{standard}} - \text{baseline}_{\text{atlantis}}} \quad (3)$$

848 Note that normalized improvement can slightly exceed 1.0 if, by chance, `Atlantis` cities perform better than the average
 849 pretrained city on some task.
 850

851 C.2. Representation Extraction

852 We extract representations from the residual stream after transformer blocks, specifically at layers 3, 4, 5, and 6 of our
 853 6-layer model. Unless otherwise specified, all representation analyses in this paper use layer 5 representations.
 854

855 To extract city representations, we pass a task prefix followed by a city ID through the model. For single-task models, we
 856 use the corresponding task prefix. For multi-task models (2-task and 3-task), we use the first task in the combination as the
 857 prefix. We verified that the choice of task prefix has negligible effect on the extracted city representations.
 858

859 For a city with ID 1234, the input sequence is:
 860

861 <bos> d i s t (c - 1 2 3 4 ,

862 We extract and concatenate the representations of two tokens: (1) the last digit of the city ID and (2) the
 863 following delimiter token (typically a comma). This yields a 256-dimensional representation (128×2) per city, which we
 864 use for both PCA visualization and linear probing.
 865

866 **Omitting cities with leading zeros** We omit cities with IDs starting with 0, 00, or 000 from representation analyses.
 867 These cities form distinct clusters in representation space, separate from cities with IDs starting with non-zero digits. We
 868 hypothesize this occurs because the digit 0 has special semantic status: in numerical outputs (distances, angles, areas),
 869 leading zeros never appear (e.g., “=769” not “=0769”), so the model learns to treat 0 differently when it appears as a
 870 leading digit. When 0 appears at the start of a city ID, the model may encode a feature indicating “this is an identifier, not a
 871 number,” causing these cities to cluster separately. To ensure consistent evaluation across all cities, we exclude IDs matching
 872 the pattern `^ [0] [0-9] * $` (i.e., any ID starting with zero).
 873

880 **C.3. Linear Probing & PCA**

881 We use the representations described in Sec. C.2 for both PCA visualization and linear probing.

883 **Linear Probing** We train linear probes to predict city coordinates (x, y) from the 256-dimensional representations. We
884 use a train/test split of 3250/1250 cities, training separate probes for x and y coordinates via ordinary least squares (OLS)
885 without regularization. We report R^2 scores and mean absolute error in scaled coordinate units.

887 **PCA** For visualization, we apply PCA to the representations and plot the first two or three principal components. We use
888 consistent color coding based on geographic region to enable visual comparison across models and seeds.

890 **Reconstruction Error** To quantify how well new entities (Atlantis cities) are integrated into the learned manifold,
891 we train linear probes exclusively on non-Atlantis cities and evaluate reconstruction error on held-out Atlantis
892 representations. Reconstruction error is measured as the absolute Euclidean distance between predicted and true coordinates.
893 Large reconstruction errors indicate that new entities are encoded in different subspaces than the original cities.

895 **C.4. Centered Kernel Alignment**

896 We use Centered Kernel Alignment (CKA) (Kornblith et al., 2019) to measure representational similarity between models.
897 Given two representation matrices $X \in \mathbb{R}^{n \times d_1}$ and $Y \in \mathbb{R}^{n \times d_2}$ (same n cities, potentially different dimensions), we
898 compute linear kernel matrices $K = XX^T$ and $L = YY^T$, center them, and compute:

$$901 \quad \text{CKA}(X, Y) = \frac{\langle K, L \rangle_F}{\|K\|_F \|L\|_F} \quad (4)$$

902 where $\langle \cdot, \cdot \rangle_F$ denotes the Frobenius inner product. CKA yields a similarity score in $[0, 1]$ that is invariant to orthogonal
903 transformations and isotropic scaling.

904 For each pair of models, we extract city representations (Sec. C.2) and compute CKA between the resulting matrices. We
905 filter cities to exclude Atlantis and IDs starting with zeros. We report CKA values at layers 3, 4, 5, and 6, with layer 5 as
906 the default unless otherwise specified.

907 **D. Additional Results**

908 **D.1. Qualitative Representations**

909 Fig. 8 shows PCA projections of city representations for single-task models across three random seeds (rows). The
910 distance task consistently produces characteristic thread-like structures. Angle and perimeter often form larger
911 2D manifold-like structures. triangle area tends to produce arc-shaped geometries. Compass forms local clusters
912 corresponding to directional categories, while inside produces a more global, diffuse structure.

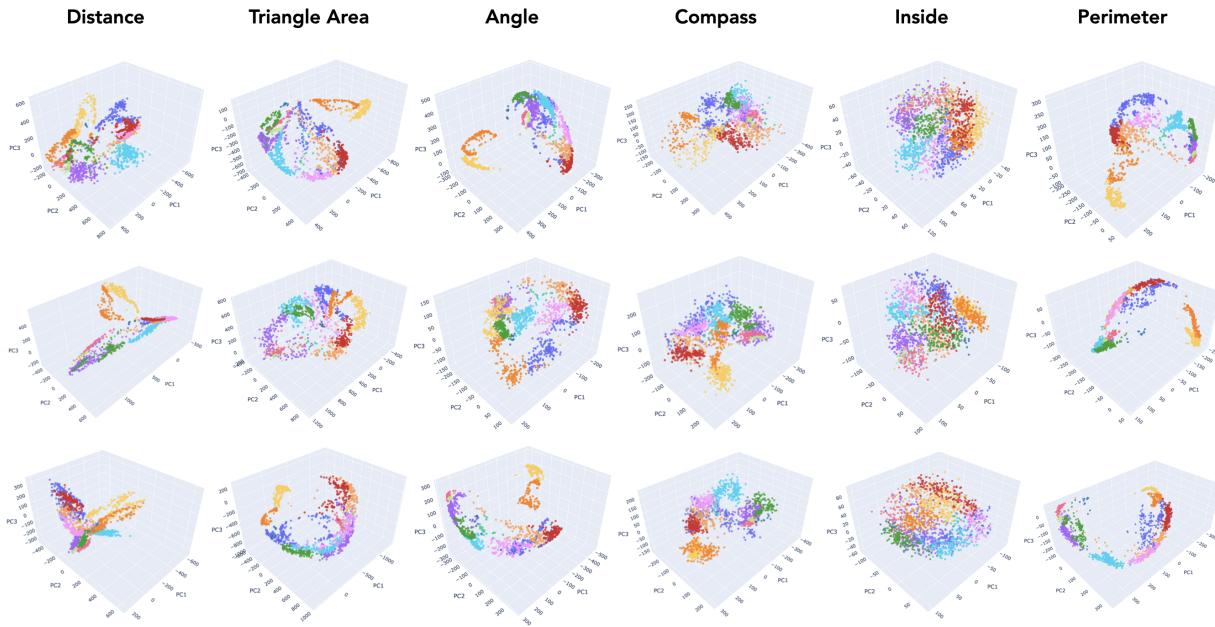
913 While there is some seed-to-seed variability within each task, the broader categories remain distinguishable: distance
914 representations are qualitatively distinct from the cluster-based representations of compass and inside, and both differ
915 from the manifold-like structures produced by triangle area, angle, and perimeter.

916 **D.2. Training Dynamics**

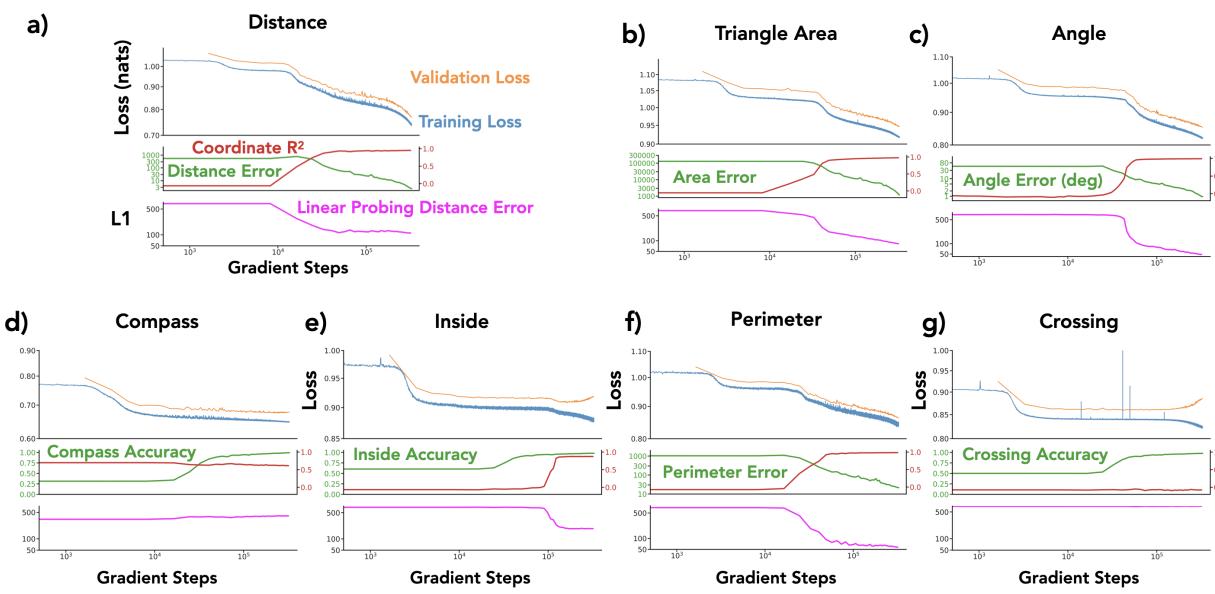
917 Fig. 9 shows training dynamics for all seven single-task models. Each panel displays three rows of metrics over gradient
918 steps: (top) training and validation loss, (middle) task-specific performance metric alongside linear probe R^2 for coordinate
919 decoding, and (bottom) linear probing distance error measuring how accurately city coordinates can be reconstructed from
920 representations.

921 Several patterns emerge across tasks. First, all tasks except crossing eventually achieve high coordinate R^2 (red curves
922 reaching ~ 1.0), indicating that world representations form reliably across diverse geometric objectives. Second, the
923 relationship between loss, task performance, and coordinate decodability varies across tasks. Third, crossing (panel g)
924 fails entirely in single-task training. Loss remains high, accuracy stays near chance, and coordinate R^2 never rises, consistent
925 with the main text observation that this task requires multi-task scaffolding.

935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957

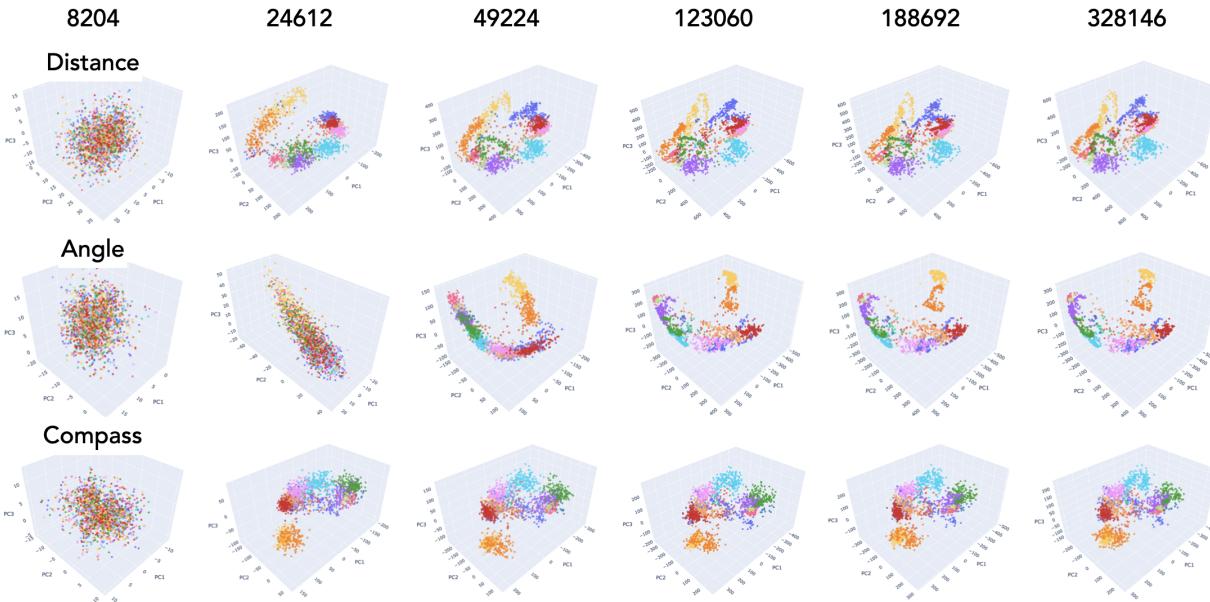


958 **Figure 8. Representation visualizations for single-task models across multiple seeds.** Each column shows a different task; each row
959 shows a different random seed. Cities are colored by geographic region.



983 **Figure 9. Training dynamics for all single-task models.** (a) distance, (b) trianglearea, (c) angle, (d) compass, (e) inside,
984 (f) perimeter, (g) crossing. Each panel shows three rows: (top) training loss (blue) and validation loss (orange); (middle)
985 task-specific metric (green, left axis) and linear probe coordinate R^2 (red, right axis); (bottom) linear probing distance error (magenta). All
986 plots use log-scale x-axis for gradient steps.
987
988
989

990
991 **Representation Dynamics.** Fig. 10 visualizes how internal representations evolve during training via PCA projections at
992 six checkpoints. A striking pattern emerges: once a representational structure forms, it remains largely fixed throughout
993 the subsequent training phase where task accuracy continues to improve. Examining the gradient steps, representations
994 are essentially fixed in the first $\sim 15\%$ of training, remaining static while loss slowly decreases and accuracy rises. The
995 distance task (top row) establishes its thread-like structure early; angle (middle row) settles into a 2D manifold;
996 compass (bottom row) forms fragmented regional clusters, all within the first few checkpoints, with minimal subsequent
997 change. What determines when representations stop evolving remains unclear, though it appears correlated with the initial
998 loss drop. This may relate to recently observed gradient dynamics in language model training, where loss deceleration
999 phases exhibit qualitatively different learning behavior (Mircea et al., 2025).



1000
1001 **Figure 10. Representation dynamics during training.** Rows: distance (top), angle (middle), compass (bottom). Columns show
1002 PCA projections at gradient steps 8204, 24612, 49224, 123060, 188692, and 328146 (left to right). Cities are colored by geographic
1003 region.
1004

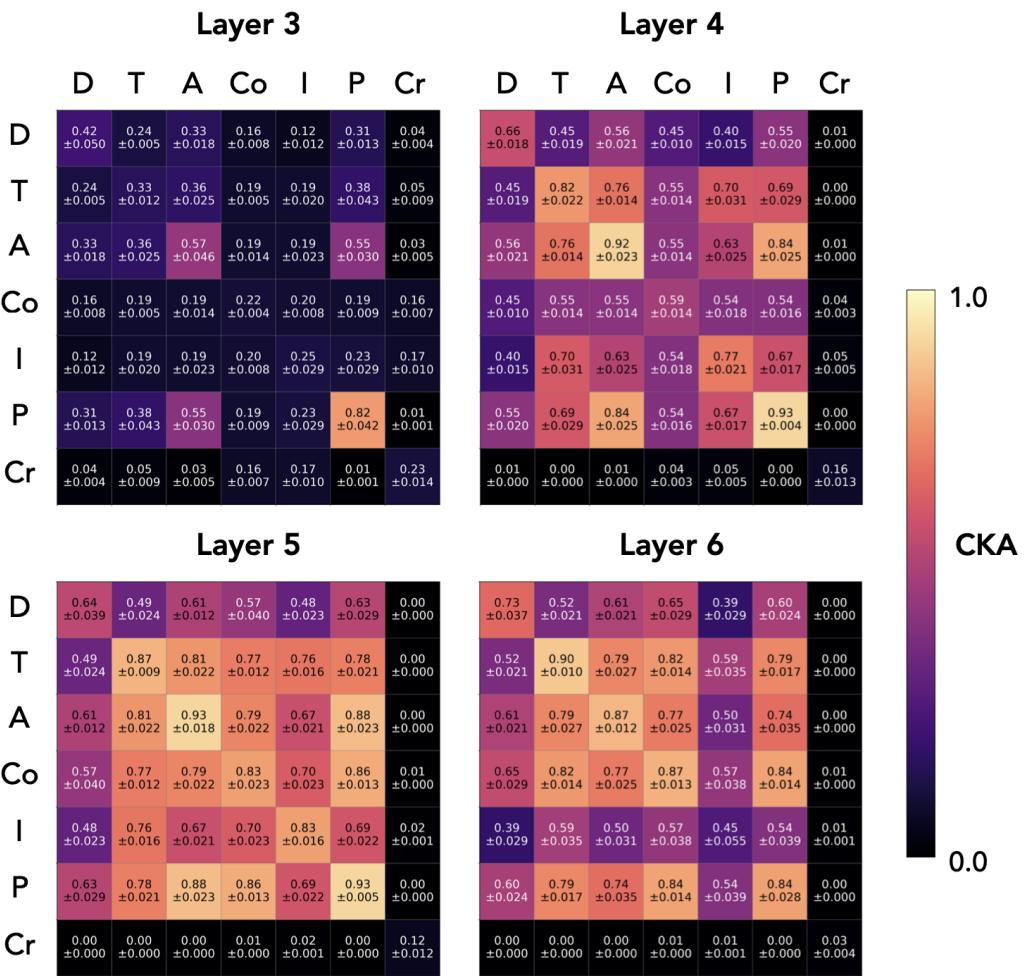
1005 D.3. Additional CKA Results

1006 **Single-Task CKA Across Layers.** Fig. 11 shows CKA matrices for single-task models at layers 3, 4, 5, and 6. Each cell
1007 shows mean \pm SEM across 3 seeds. We observe: (1) CKA values increase from layer 3 to layers 4–6, indicating that world
1008 representations become more consistent in later layers; (2) the *distance* task (D) shows lower CKA with other tasks
1009 across all layers, consistent with its divergent representational geometry; (3) *crossing* (Cr) shows near-zero CKA due to
1010 training failure in single-task settings; (4) diagonal entries (same task) can show significant variability, indicating that even
1011 identical training objectives can yield different representational solutions.
1012

1013 **Two-Task CKA.** Fig. 12 shows the CKA matrix for two-task models at layer 5. Compared to single-task models (Fig. 11,
1014 layer 5), two-task training substantially increases representational alignment: all off-diagonal entries exceed 0.84, compared
1015 to values as low as 0.48 for single-task models. Notably, diagonal entries (same task combination, different seeds) show
1016 minimum CKA of 0.89, indicating that multi-task training also reduces inter-seed variance. For diagonal entries, we exclude
1017 same-seed comparisons (which trivially yield 1.0) and report only the upper triangle since the matrix is symmetric. This
1018 confirms the main text finding that multi-task training drives representational convergence.
1019

1020 **CKA vs. Task Count (Per-Seed).** Fig. 13 shows the same CKA vs. task count analysis as Fig. 3(d) in the main text,
1021 but broken down by individual seeds. Each panel shows one seed. These per-seed values are pooled to produce the main
1022 text figure, where error bars represent SEM across seeds. The pattern is consistent across all three seeds: CKA increases
1023 substantially from 1 to 2 tasks and saturates at 2–3 tasks for layers 4–6.
1024

1045
1046
1047
1048
1049
1050
1051
1052
1053
1054
1055
1056
1057
1058
1059
1060
1061



1077 **Figure 11. CKA matrices for single-task models across layers.** Each cell shows mean \pm SEM across 3 seeds. D=distance, T=triangle area, A=angle, Co=compass, I=inside, P=perimeter, Cr=crossing. CKA increases in later layers; distance shows consistently lower cross-task similarity.

1080
1081
1082
1083
1084
1085
1086
1087
1088
1089
1090
1091
1092
1093
1094
1095
1096
1097
1098
1099

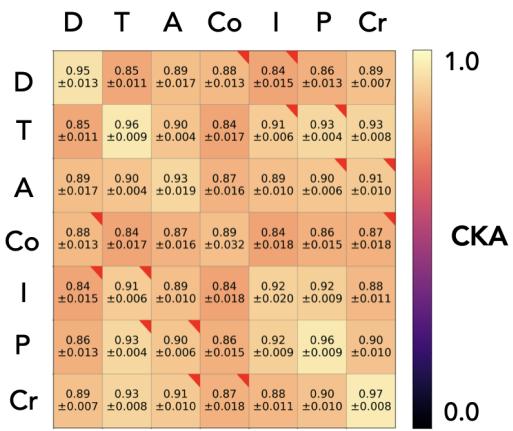


Figure 12. CKA matrix for two-task models at layer 5. Mean \pm SEM across 3 seeds. All pairs show high alignment (>0.84), substantially higher than single-task models.

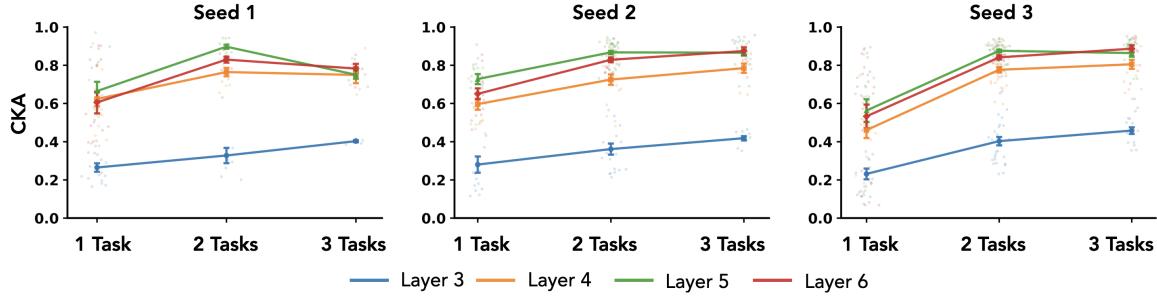


Figure 13. CKA vs. task count for individual seeds. Each panel shows a different seed. These values are pooled in Fig. 3(d); error bars there represent SEM across seeds.

Aggregated CKA Trends. Fig. 14(a) shows CKA vs. task count for a single seed, using all $\binom{7}{2} = 21$ two-task models and all $\binom{7}{3} = 35$ three-task models, but only comparing non-overlapping pairs (models sharing no common tasks). This yields 105 non-overlapping pairs for 2-task models and 70 for 3-task models. Fig. 14(b) shows within-task CKA (same task combination, different seeds) as a function of task count, demonstrating that multi-task training also reduces seed-to-seed variability: representations become more consistent not just across tasks but also across random initializations.

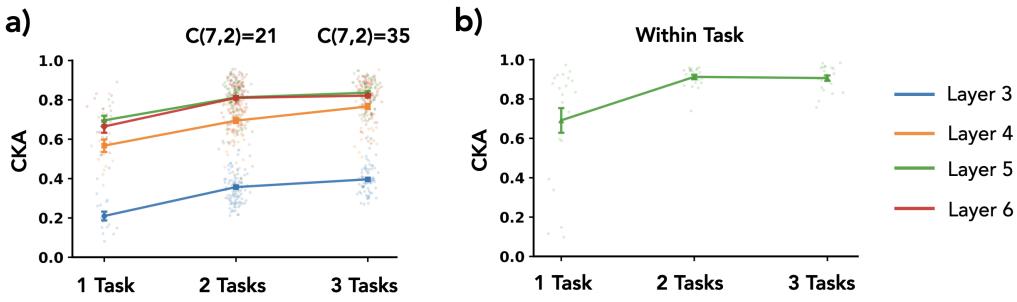


Figure 14. Aggregated CKA analysis. (a) CKA vs. task count for single seed, comparing only non-overlapping model pairs (105 pairs for 2-task, 70 pairs for 3-task). (b) Within-task CKA (same task combination, different seeds) increases with task count, indicating multi-task training reduces seed variability.

CKA vs. Generalization (Annotated). Fig. 15 is an annotated version of Fig. 5(b), with each point labeled by its (train→eval) task pair.

D.4. Additional Fine-Tuning Evaluation Results

Raw fine-tuning results for individual seeds.

D.5. Pretraining Variations

Pretraining with Atlantis. In the main text, we showed that fine-tuning on divergent tasks fails to integrate Atlantis cities into the learned representation manifold (Fig. 6d, red histogram). To verify that this failure stems from fine-tuning dynamics rather than a peculiarity of the geometry around Atlantis, we trained a model with Atlantis cities included from the start of pretraining. Fig. 19 shows the resulting representations: Atlantis cities are seamlessly integrated into the world manifold, indistinguishable from other cities in both PCA projections (a) and linear probe reconstructions (b). This confirms that the representation space can readily accommodate Atlantis, and thus, the integration failure observed in fine-tuning is a property of the optimization dynamics, not a fundamental limitation of the architecture or task.

Wider Model. To test whether our findings depend on model capacity, we trained a wider model with $2\times$ the hidden dimension (256 vs. 128) and intermediate size (1024 vs. 512), resulting in approximately $4\times$ the parameters. Fig. 20 shows fine-tuning results for this wider model: (a) single-task fine-tuning normalized improvement; (b) two-task fine-

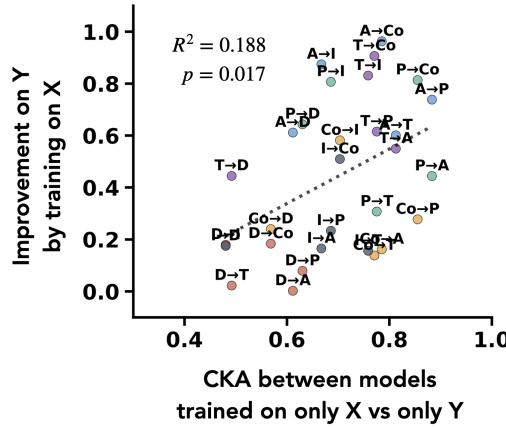


Figure 15. Annotated version of Fig. 5(b). Each point is labeled with its (train \rightarrow eval) task pair. D=distance, T=triangle area, A=angle, Co=compass, I=inside, P=perimeter.

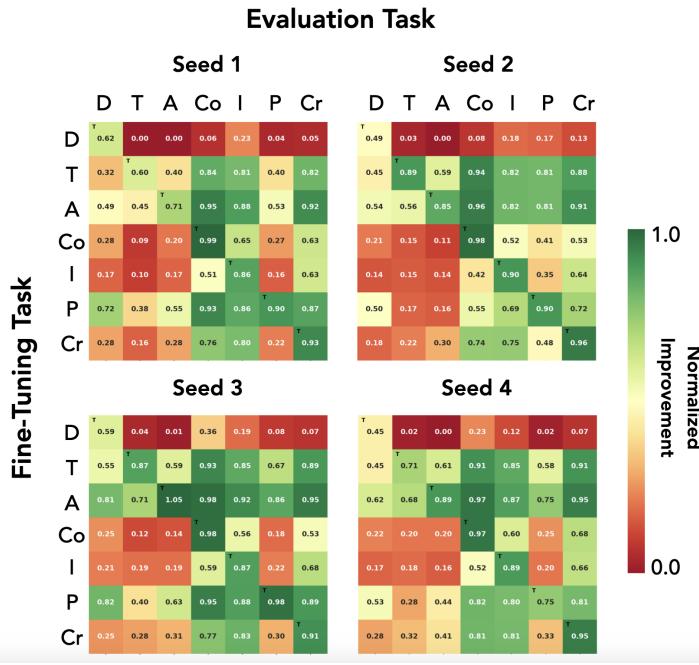


Figure 16. Single-task fine-tuning results for individual seeds. Per-seed version of Fig. 5(a), organized in a 2×2 grid.

1210
1211
1212
1213
1214
1215
1216
1217
1218
1219
1220
1221
1222
1223
1224
1225
1226
1227
1228
1229
1230



Figure 17. Two-task fine-tuning normalized improvement for all 21 task combinations. Leftmost panel shows average across seeds; remaining panels show individual seeds.

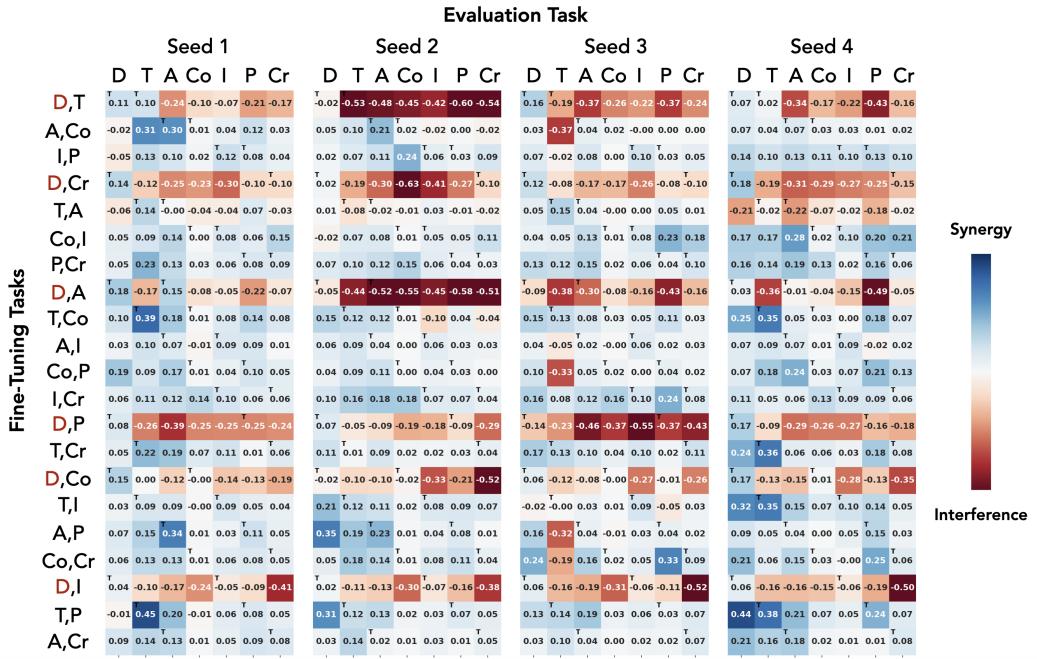


Figure 18. Deviation from best-teacher expectation for all 21 two-task combinations. All 4 seeds shown; average is in main text Fig. 6(c).

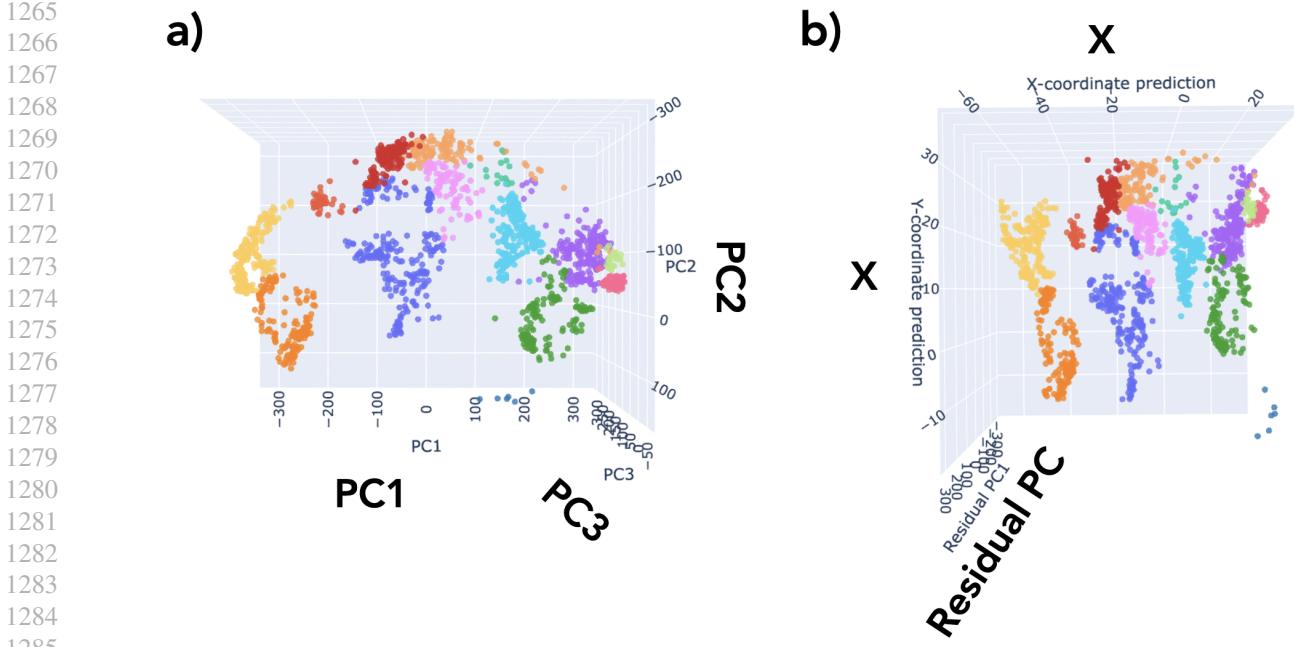


Figure 19. Representations when **Atlantis** is included during pretraining. (a) PCA projection showing **Atlantis** cities (small cluster in Atlantic region) integrated with world cities. (b) Linear probe reconstruction confirming geographic accuracy. Unlike fine-tuned models, **Atlantis** cities lie on the same manifold as other cities.

tuning normalized improvement; (c) deviation from best-teacher expectation. We still observe that distance-containing combinations (red labels in panel c) show degraded cross-task generalization. This suggests that divergent task interference is not simply a capacity limitation.

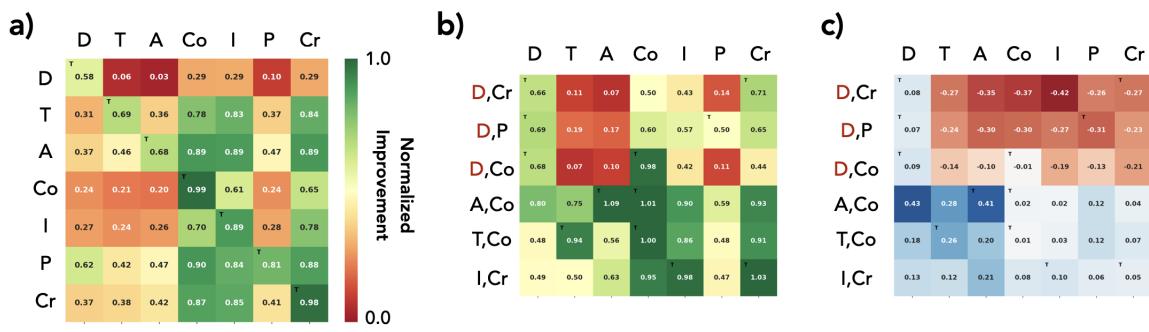


Figure 20. Fine-tuning results for wider model (2× hidden dimension). For all panels: rows = fine-tuning task(s), columns = evaluation task. (a) Single-task fine-tuning normalized improvement. (b) Two-task fine-tuning normalized improvement. (c) Deviation from best-teacher expectation; distance-containing combinations (red labels) still show degraded generalization.