

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 Transformer 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 and all cities, then test whether new entities can be consistently integrated into the representation space via fine-tuning. Surprisingly, we find that despite multi-task pretraining, some tasks—which we dub *divergent*—actively harm the representational integration of new entities. Our results show that training on multiple relational tasks reliably produces convergent world representations, but some 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;

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.

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. Even the most accessible factor, the model, becomes expensive to perturb at scale; the data is harder still to control; and the underlying world is effectively immutable. 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:

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

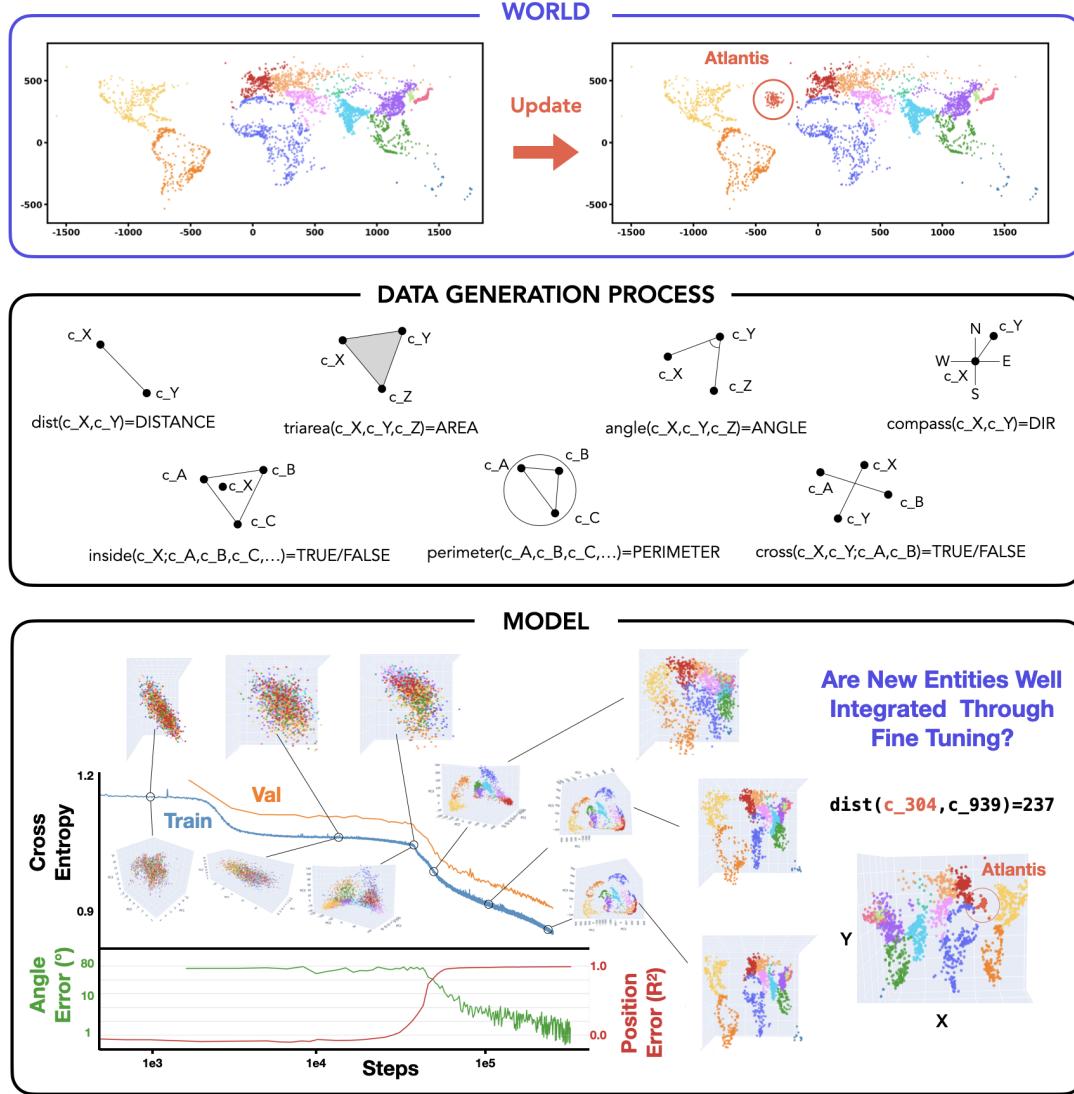


Figure 1. Overview of the World-Data-Model framework. **Top:** The world consists of 5,075 real city coordinates; we test adaptation by adding 100 synthetic Atlantis cities (App. A.1). **Middle:** Seven geometric tasks generate training data from city coordinates (App. A.2). **Bottom:** Training dynamics 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.

- **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 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.

2. Related Work

Internal Representations. Understanding internal representations has been fundamental since the development of neural networks (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). Mechanistic interpretability and sparse autoencoders have further enabled decomposition of neural activations into interpretable features (Anthropic AI, 2023; Templeton et al., 2024). 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). 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). 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). 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; Raventós 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.

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:

- Learnability:** All tasks are deterministically generated from the same underlying coordinates. A model that learns the world structure can leverage it across all tasks.
- Latent State:** Models never see coordinates directly, only task outputs, allowing us to probe whether they internally reconstruct the world structure.
- 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

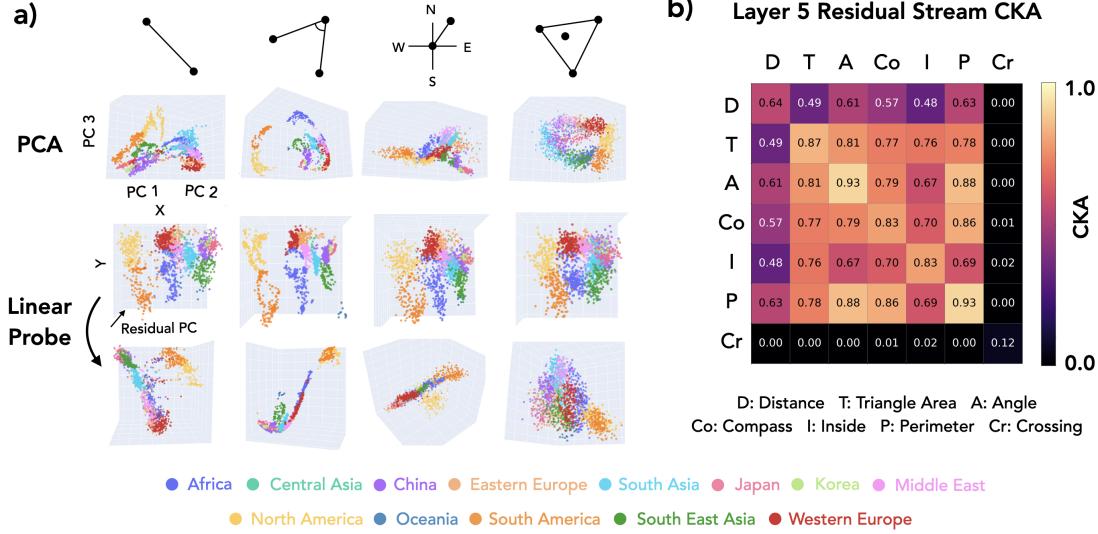


Figure 2. World representation geometry depends on the data generation process. (a) Different tasks create distinct geometries: distance (thread-like), angle (2D manifold), compass (fragmented), inside (diffuse). Row 1: PCA. Row 2: Linear probe projections. Row 3: Rotated views showing hidden structure. 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.

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 $\text{cross}(c_2345, c_6789; c_0123, c_4567) = \text{TRUE}$ 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. A.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

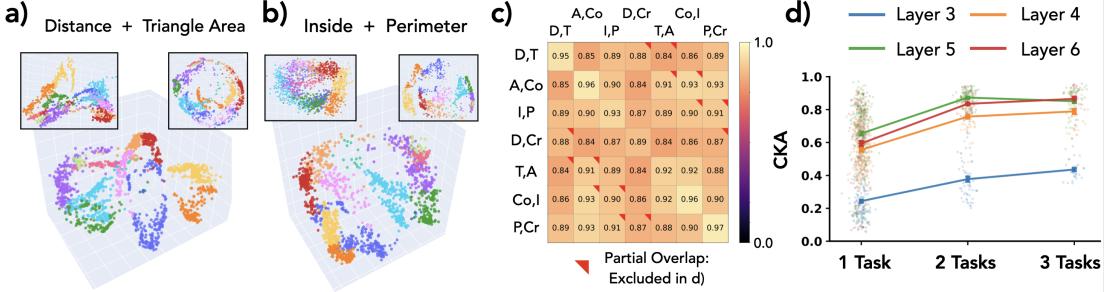


Figure 3. Multi-task pretraining drives representational convergence. (a,b) Two-task training creates more regular structures than single-task models. (c) CKA matrix (7×7) for two-task models shows higher alignment (see App. Fig. 12 for SEM). (d) Average CKA increases with task count ($1 \rightarrow 2 \rightarrow 3$), 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.

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. C.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 representations 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).

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

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).

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 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.

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

275 (App. Fig. 14b).

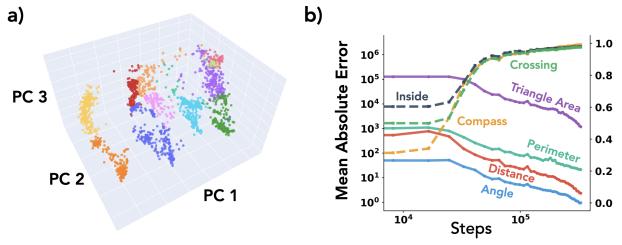


Figure 4. 7-task model. (a) PCA projection of layer 5 representations naturally reveals world map structure. (b) Training curves showing successful learning of all 7 tasks, including crossing which failed in single-task training.

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

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

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).

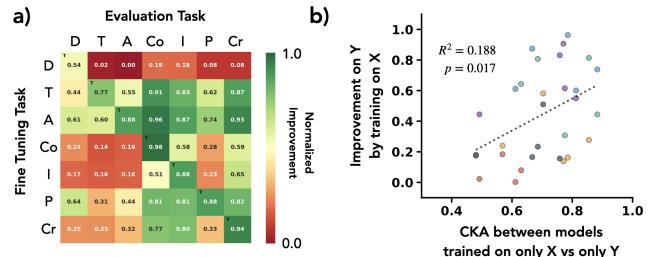


Figure 5. Fine-tuning generalization and its correlation with representational similarity. (a) Generalization matrix (averaged over 4 seeds; see App. Fig. 16 for individual seeds): each row is a model that integrated Atlantis via one task; columns show normalized improvement on Atlantis queries for each task (see App. B.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. A.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:

"Even though models develop joint world representations which converge in multi-task pretraining, gradient descent on divergent tasks might fail to act on these shared representations during fine-tuning, instead utilizing hidden spaces that don't propagate updates across tasks."

Our question is then two-part:

1. To what extent do divergent tasks affect fine-tuning generalization?
2. Will gradient descent on divergent tasks fail to merge fine-tuning introduced concepts to the original representation manifold?

Result 2: Divergent Tasks Catastrophically Harm Generalization To investigate how divergent tasks affect generalization, we move from single-task to multi-task fine-tuning settings. First, we introduce a simple heuristic model: fine-tuning on a concatenated dataset $\{D_1, D_2, \dots, D_n\}$ (which do not provide conflicting supervision) should combine their individual effects. Specifically, when concatenating and shuffling all fine-tuning data to avoid sequential learning effects like catastrophic forgetting (McCloskey & Cohen, 1989), we expect the improvement Imp_i on task i after training on tasks j and k to follow a **best-teacher model**:

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

To test this hypothesis, we fine-tuned the 7-task model on all $\binom{7}{2} = 21$ possible two-task combinations. Fig. 6(a,c) shows the *deviation* from our best-teacher expectation (averaged over 4 seeds; see App. Fig. 17 for raw improvements and App. Fig. 18 for individual seeds). Strikingly, we observe "red horizontal bands", models that not only fail to benefit from multi-task training but actually perform worse

than their best single-task component. Notably, all these degraded performance bands involve the *distance* task. Fig. 6(c) quantifies this: when we split the deviation values into models with and without *distance*, we consistently observe lower-than-expected performance when the divergent task is included. This confirms that *divergent tasks (those with low single-task CKA) actively harm fine-tuning generalization rather than simply failing to contribute*. We next examine how this manifests in the learned representations.

Result 3: Divergent Tasks Disrupt Representational Integration of New Entities Having shown that divergent tasks harm generalization (Question 1), we now address Question 2: does gradient descent on divergent tasks fail to merge new entities into the representation manifold?

We take two exemplars from the 21 fine-tuning runs: one without *distance* that generalized well (*angle + compass*), and one with *distance* that was harmed (*distance + perimeter*). We first train a linear probe on a subset of all cities including *Atlantis*; these reconstructions are shown in Fig. 6(b) (top and bottom panels). In the well-integrated case, *Atlantis* cities lie within the world data manifold. In the ill-integrated case, *Atlantis* cities are off the manifold. While this difference appears subtle in 2D projections, the effect is dramatic in 3D, we strongly encourage readers to explore our interactive visualizations. Next, we train a linear probe on 4000 non-*Atlantis* cities and apply it to *Atlantis* representations (middle panels). In the well-integrated case, *Atlantis* cities (red-orange) are relatively well reconstructed compared to ground truth (black crosses); in the ill-integrated case, reconstruction fails completely.

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

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

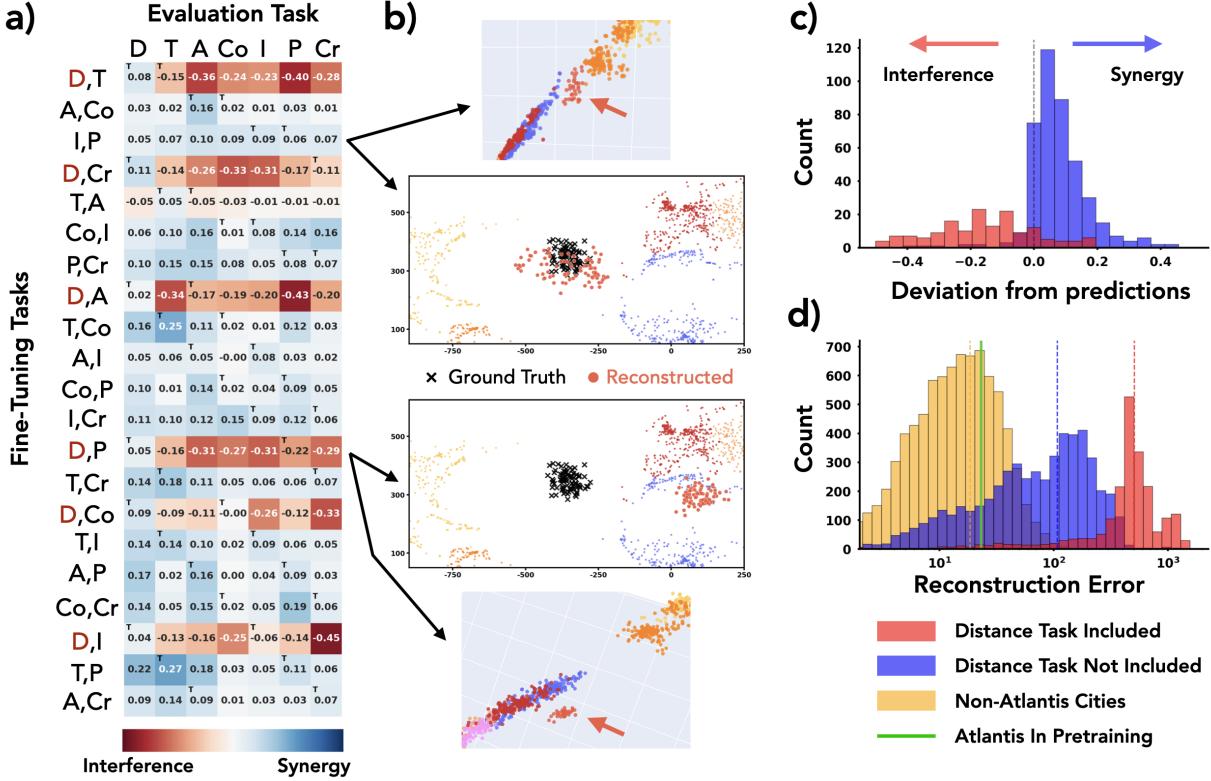


Figure 6. Divergent tasks harm multi-task fine-tuning and disrupt representational integration. (a) Deviation from best-teacher expectation for 21 two-task models (rows) across 7 evaluation tasks (columns), computed in normalized improvement space (see App. B.1); “red horizontal bands” show distance 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.

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

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

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., 2025). 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., 2025; 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*

440
 441 *the forward pass does not imply modularity in the back-*
 442 *ward pass.* Multi-task training produces clean, structured
 443 representations that can be easily decoded into world coor-
 444 dinates, yet these world models can be fractured and partial
 445 when it comes to adaptation. Gradient descent may not
 446 respect the forward-pass modularity when updating weights:
 447 fine-tuning on divergent tasks routes updates through path-
 448 ways that bypass the shared world manifold, encoding new
 449 entities in task-specific subspaces.

450
 451 **Future work.** Understanding the mechanistic basis of task
 452 divergence is an important open question. If divergence is
 453 a property of task-architecture pairing rather than learned
 454 weights, it may be predictable from task structure and gra-
 455 dient geometry alone, enabling identification of harmful tasks
 456 before training.

457 **Limitations.** We study world representation formation and
 458 adaptation in a controlled synthetic setting with small-scale
 459 models. While we find non-trivial phenomenology, includ-
 460 ing the emergence of world representations, task-dependent
 461 geometry, representational convergence under multi-task
 462 training, and off-target fine-tuning effects, it is difficult
 463 to guarantee these findings will generalize to large-scale
 464 models trained on natural data. Additionally, we identify
 465 divergence as a diagnostic marker but do not reveal under-
 466 lying mechanisms. Furthermore, our claims regarding the
 467 Platonic Representation Hypothesis are partial: we demon-
 468 strate task-driven convergence within a single architecture
 469 and modality, but do not explore true multimodality or cross-
 470 architecture convergence.

471 7. Conclusion

472 We introduced a framework separating world from data
 473 generation to study how representations form and adapt. In-
 474 creasing task diversity drives **convergent** world representa-
 475 tions: models trained on multiple tasks develop increasingly
 476 aligned geometry, even when sharing no common tasks,
 477 supporting the Multitask Scaling Hypothesis of PRH. Yet
 478 convergence masks vulnerability: certain **divergent** tasks
 479 actively harm integration of new entities during fine-tuning,
 480 encoding them in hidden spaces rather than the shared man-
 481 ifold. Clean representations do not guarantee clean adapta-
 482 tion.

483 Impact Statement

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

References

- 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>.
- 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>.
- 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>.
- 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>.

- 495 Chan, S. C. Y., Santoro, A., Lampinen, A. K., Wang, J. X.,
 496 Singh, A., Richemond, P. H., McClelland, J., and Hill, F.
 497 Data distributional properties drive emergent in-context
 498 learning in transformers, 2022. URL <https://arxiv.org/abs/2205.05055>.
- 500 Devlin, J., Chang, M.-W., Lee, K., and Toutanova, K. Bert:
 501 Pre-training of deep bidirectional transformers for lan-
 502 guage understanding. *arXiv preprint arXiv:1810.04805*,
 503 2018.
- 504 Dohare, S., Hernandez-Garcia, J. F., Rahman, P., Mah-
 505 moud, A. R., and Sutton, R. S. Maintaining plastic-
 506 ity in deep continual learning, 2024. URL <https://arxiv.org/abs/2306.13812>.
- 507 Fukushima, K. Neocognitron: A self-organizing neural
 508 network model for a mechanism of pattern recognition
 509 unaffected by shift in position. *Biological cybernetics*, 36
 510 (4):193–202, 1980.
- 511 Gopalani, P. and Hu, W. What happens during the loss
 512 plateau? understanding abrupt learning in transfor-
 513 mers, 2025. URL <https://arxiv.org/abs/2506.13688>.
- 514 Gurnee, W. and Tegmark, M. Language models represent
 515 space and time. *arXiv preprint arXiv:2310.02207*, 2023.
- 516 He, K., Zhang, X., Ren, S., and Sun, J. Deep residual
 517 learning for image recognition, 2015.
- 518 Hindupur, S. S. R., Lubana, E. S., Fel, T., and Ba, D.
 519 Projecting assumptions: The duality between sparse au-
 520 toencoders and concept geometry, 2025. URL <https://arxiv.org/abs/2503.01822>.
- 521 Hoffmann, D. T., Schrodi, S., Bratulić, J., Behrmann, N.,
 522 Fischer, V., and Brox, T. Eureka-moments in transfor-
 523 mers: Multi-step tasks reveal softmax induced optimiza-
 524 tion problems, 2024. URL <https://arxiv.org/abs/2310.12956>.
- 525 Hubel, D. H. and Wiesel, T. N. Receptive fields, binocular
 526 interaction and functional architecture in the cat’s visual
 527 cortex. *The Journal of physiology*, 160(1):106, 1962.
- 528 Huh, M., Cheung, B., Wang, T., and Isola, P. The pla-
 529 tonic representation hypothesis, 2024. URL <https://arxiv.org/abs/2405.07987>.
- 530 Jain, S., Kirk, R., Lubana, E. S., Dick, R. P., Tanaka,
 531 H., Grefenstette, E., Rocktäschel, T., and Krueger,
 532 D. S. Mechanistically analyzing the effects of fine-
 533 tuning on procedurally defined tasks. *arXiv preprint
 534 arXiv:2311.12786*, 2023.
- 535 Kornblith, S., Norouzi, M., Lee, H., and Hinton, G. Similar-
 536 ity of Neural Network Representations Revisited. In *Proc.
 537 of the 36th Proc. Int. Conf. on Machine Learning (ICML)*,
 538 Proc. of Machine Learning Research. PMLR, 09–15 Jun
 539 2019.
- 540 Krizhevsky, A., Sutskever, I., and Hinton, G. E. Imagenet
 541 classification with deep convolutional neural networks.
 542 *Advances in neural information processing systems*, 25,
 543 2012.
- 544 Kumar, A., Raghunathan, A., Jones, R., Ma, T., and Liang,
 545 P. Fine-tuning can distort pretrained features and un-
 546 derperform out-of-distribution, 2022. URL <https://arxiv.org/abs/2202.10054>.
- 547 Kumar, A., Clune, J., Lehman, J., and Stanley, K. O. Ques-
 548 tioning representational optimism in deep learning: The
 549 fractured entangled representation hypothesis, 2025. URL <https://arxiv.org/abs/2505.11581>.
- 550 Lampinen, A. K., Chaudhry, A., Chan, S. C. Y., Wild, C.,
 551 Wan, D., Ku, A., Bornschein, J., Pascanu, R., Shanahan,
 552 M., and McClelland, J. L. On the generalization
 553 of language models from in-context learning and
 554 finetuning: a controlled study, 2025. URL <https://arxiv.org/abs/2505.00661>.
- 555 Li, K., Hopkins, A. K., Bau, D., Viégas, F., Pfister, H.,
 556 and Wattenberg, M. Emergent world representations:
 557 Exploring a sequence model trained on a synthetic task.
 558 In *The Eleventh International Conference on Learning
 559 Representations*, 2022.
- 560 Li, Y., Campbell, D., Chan, S. C. Y., and Lampinen, A. K.
 561 Just-in-time and distributed task representations in lan-
 562 guage models, 2025. URL <https://arxiv.org/abs/2509.04466>.
- 563 Lubana, E. S., Rager, C., Hindupur, S. S. R., Costa, V.,
 564 Tuckute, G., Patel, O., Murthy, S. K., Fel, T., Wurgafit, D.,
 565 Bigelow, E. J., Lin, J., Ba, D., Wattenberg, M., Viegas,
 566 F., Weber, M., and Mueller, A. Priors in time: Missing
 567 inductive biases for language model interpretability, 2025.
 568 URL <https://arxiv.org/abs/2511.01836>.
- 569 Marks, S. and Tegmark, M. The geometry of truth: Emer-
 570 gent linear structure in large language model represen-
 571 tations of true/false datasets, 2024. URL <https://arxiv.org/abs/2310.06824>.
- 572 McCloskey, M. and Cohen, N. J. Catastrophic interfer-
 573 ence in connectionist networks: The sequential learning
 574 problem. In *Psychology of learning and motivation*, vol-
 575 ume 24, pp. 109–165. Elsevier, 1989.
- 576 Menon, A., Shrivastava, M., Krueger, D., and Lubana, E. S.
 577 Analyzing (in)abilities of saes via formal languages, 2025.
 578 URL <https://arxiv.org/abs/2410.11767>.

- 550 Michaud, E. J., Liu, Z., Girit, U., and Tegmark, M. The
551 quantization model of neural scaling. *arXiv preprint*
552 *arXiv:2303.13506*, 2023.
- 553 Minder, J., Dumas, C., Juang, C., Chugtai, B., and Nanda,
554 N. Overcoming sparsity artifacts in crosscoders to inter-
555 pret chat-tuning, 2025. URL <https://arxiv.org/abs/2504.02922>.
- 556 Nanda, N., Chan, L., Lieberum, T., Smith, J., and Stein-
557 hardt, J. Progress measures for grokking via mechanistic
558 interpretability, 2023a. URL <https://arxiv.org/abs/2301.05217>.
- 559 Nanda, N., Lee, A., and Wattenberg, M. Emergent lin-
560 ear representations in world models of self-supervised
561 sequence models. In *Proceedings of the 6th Black-*
562 *boxNLP Workshop: Analyzing and Interpreting Neural*
563 *Networks for NLP*, pp. 16–30, 2023b. URL <https://arxiv.org/abs/2309.00941>.
- 564 Nishi, K., Okawa, M., Ramesh, R., Khona, M., Lubana,
565 E. S., and Tanaka, H. Representation shattering in trans-
566 formers: A synthetic study with knowledge editing. *arXiv*
567 *preprint arXiv:2410.17194*, 2024.
- 568 Okawa, M., Lubana, E. S., Dick, R. P., and Tanaka, H. Com-
569 positional abilities emerge multiplicatively: Exploring
570 diffusion models on a synthetic task, 2024.
- 571 Park, C. F., Lee, A., Lubana, E. S., Yang, Y., Okawa, M.,
572 Nishi, K., Wattenberg, M., and Tanaka, H. Iclr: In-context
573 learning of representations, 2024a. URL <https://arxiv.org/abs/2501.00070>.
- 574 Park, C. F., Lubana, E. S., Pres, I., and Tanaka, H. Com-
575 petition dynamics shape algorithmic phases of in-context
576 learning. *arXiv preprint arXiv:2412.01003*, 2024b.
- 577 Park, C. F., Okawa, M., Lee, A., Lubana, E. S., and Tanaka,
578 H. Emergence of hidden capabilities: Exploring learning
579 dynamics in concept space, 2024c. URL <https://arxiv.org/abs/2406.19370>.
- 580 Park, C. F., Zhang, Z., and Tanaka, H. *New News*: System-2
581 fine-tuning for robust integration of new knowledge, 2025.
582 URL <https://arxiv.org/abs/2505.01812>.
- 583 Pezeshki, M., Kaba, O., Bengio, Y., Courville, A. C., Pre-
584 cup, D., and Lajoie, G. Gradient starvation: A learning
585 proclivity in neural networks. *Adv. in Neural Information*
586 *Processing Systems (NeurIPS)*, 2021.
- 587 Qin, T., Park, C. F., Kwun, M., Walsman, A., Malach, E.,
588 Anand, N., Tanaka, H., and Alvarez-Melis, D. Decom-
589 posing elements of problem solving: What "math" does
590 rl teach?, 2025. URL <https://arxiv.org/abs/2505.22756>.
- 591 Radford, A., Narasimhan, K., Salimans, T., Sutskever, I.,
592 et al. Improving language understanding by generative
593 pre-training, 2018.
- 594 Raventós, A., Paul, M., Chen, F., and Ganguli, S. Pretrain-
595 ing task diversity and the emergence of non-bayesian
596 in-context learning for regression, 2023. URL <https://arxiv.org/abs/2306.15063>.
- 597 Reddy, G. The mechanistic basis of data dependence and
598 abrupt learning in an in-context classification task, 2023.
599 URL <https://arxiv.org/abs/2312.03002>.
- 600 Rosenblatt, F. The perceptron: a probabilistic model for
601 information storage and organization in the brain. *Psy-
602 chological review*, 65(6):386, 1958.
- 603 Rumelhart, D. E., Hinton, G. E., and Williams, R. J. Learn-
604 ing representations by back-propagating errors. *nature*,
605 323(6088):533–536, 1986.
- 606 Shah, H., Tamuly, K., Raghunathan, A., Jain, P., and Ne-
607 trapalli, P. The pitfalls of simplicity bias in neural
608 networks, 2020. URL <https://arxiv.org/abs/2006.07710>.
- 609 Shai, A. S., Marzen, S. E., Teixeira, L., Oldenziel, A. G.,
610 and Riechers, P. M. Transformers represent belief state
611 geometry in their residual stream, 2025. URL <https://arxiv.org/abs/2405.15943>.
- 612 Templeton, A., Conerly, T., Marcus, J., Lindsey, J., Bricken,
613 T., Chen, B., Pearce, A., Citro, C., Ameisen, E., Jones,
614 A., Cunningham, H., Turner, N. L., McDougall, C.,
615 MacDiarmid, M., Freeman, C. D., Sumers, T. R.,
616 Rees, E., Batson, J., Jermyn, A., Carter, S., Olah,
617 C., and Henighan, T. Scaling monosemanticity: Ex-
618 tracting interpretable features from claude 3 sonnet.
619 *Transformer Circuits Thread*, 2024. URL <https://transformer-circuits.pub/2024/scaling-monosemanticity/index.html>.
- 620 Wang, A., Engels, J., Clive-Griffin, O., Rajamanoharan, S.,
621 and Nanda, N. Simple mechanistic explanations for out-
622 of-context reasoning, 2025. URL <https://arxiv.org/abs/2507.08218>.
- 623 Ward, J., Lin, C., Venhoff, C., and Nanda, N. Reasoning-
624 finetuning repurposes latent representations in base mod-
625 els, 2025. URL <https://arxiv.org/abs/2507.12638>.
- 626 Wurgaft, D., Lubana, E. S., Park, C. F., Tanaka, H., Reddy,
627 G., and Goodman, N. D. In-context learning strategies
628 emerge rationally, 2025. URL <https://arxiv.org/abs/2506.17859>.

605 Xie, S. M., Raghunathan, A., Liang, P., and Ma, T. An
606 explanation of in-context learning as implicit bayesian
607 inference. *arXiv preprint arXiv:2111.02080*, 2021.

608
609 Yue, Y., Chen, Z., Lu, R., Zhao, A., Wang, Z., Yue, Y.,
610 Song, S., and Huang, G. Does reinforcement learn-
611 ing really incentivize reasoning capacity in llms beyond
612 the base model?, 2025. URL <https://arxiv.org/abs/2504.13837>.

613
614 Zhang, S., Patel, A., Rizvi, S. A., Liu, N., He, S., Karbasi,
615 A., Zappala, E., and van Dijk, D. Intelligence at the edge
616 of chaos, 2025. URL <https://arxiv.org/abs/2410.02536>.

618

619

620

621

622

623

624

625

626

627

628

629

630

631

632

633

634

635

636

637

638

639

640

641

642

643

644

645

646

647

648

649

650

651

652

653

654

655

656

657

658

659

APPENDIX

A. Experimental Details

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

A.1. World

Our experiments use a geographic world consisting of 5,075 cities extracted from the GeoNames database 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 problems.

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

A.2. Data Generation Process

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

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

A.3. Model and Training

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

Task	Input	Output Type	Unit/Values	Example
distance	2 cities	Numerical	Scaled coords	dist (c ₁ , c ₂)
triarea	3 cities	Numerical	Scaled coords ²	triarea (c ₁ , c ₂ , c ₃)
angle	3 cities	Numerical	Degrees (0–180)	angle (c ₁ , c ₂ , c ₃)
compass	2 cities	Categorical	8 directions	compass (c ₁ , c ₂)
inside	1 + n cities	Categorical	TRUE/FALSE	inside (c ₁ , c ₂ , ..., c _n)
perimeter	n cities	Numerical	Scaled coords	perimeter (c ₁ , c ₂ , ..., c _n)
crossing	4 cities	Categorical	TRUE/FALSE	crossing (c ₁ , c ₂ , c ₃ , c ₄)

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

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

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

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)

Table 2. Pretraining hyperparameters.

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

B. Analysis Methods

B.1. Evaluation

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

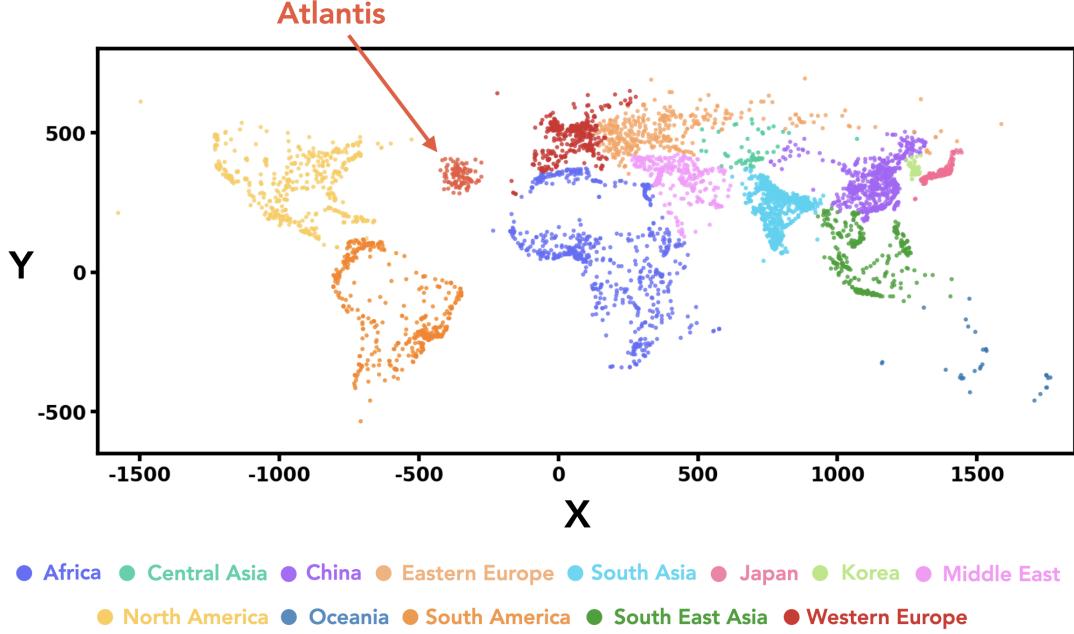


Figure 7. **Geographic distribution of cities used in our experiments.** 5,075 real-world cities plus 100 synthetic Atlantis cities (5,175 total). Cities span all continents and provide a fixed, measurable world structure. Coordinates use an equirectangular projection: $x = 10 \times \text{longitude}$,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

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

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

B.2. Representation Extraction

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

B.3. Linear Probing & PCA

We train linear probes to predict city coordinates (x, y) from the 256-dimensional representations. We use a train/test split of 3250/1250 cities, training separate probes for x and y coordinates via ordinary least squares (OLS) without regularization.

B.4. Centered Kernel Alignment

We use Centered Kernel Alignment (CKA) to measure representational similarity between models. CKA yields a similarity score in $[0, 1]$ that is invariant to orthogonal transformations and isotropic scaling.

C. Additional Results

C.1. Qualitative Representations

Fig. 8 shows PCA projections of city representations for single-task models across three random seeds.

C.2. Training Dynamics

Fig. 9 shows training dynamics for all seven single-task models.

Representation Dynamics. Fig. 10 visualizes how internal representations evolve during training via PCA projections at six checkpoints.

C.3. Additional CKA Results

Single-Task CKA Across Layers. Fig. 11 shows CKA matrices for single-task models at layers 3, 4, 5 and 6.

Two-Task CKA. Fig. 12 shows the CKA matrix for two-task models at layer 5.

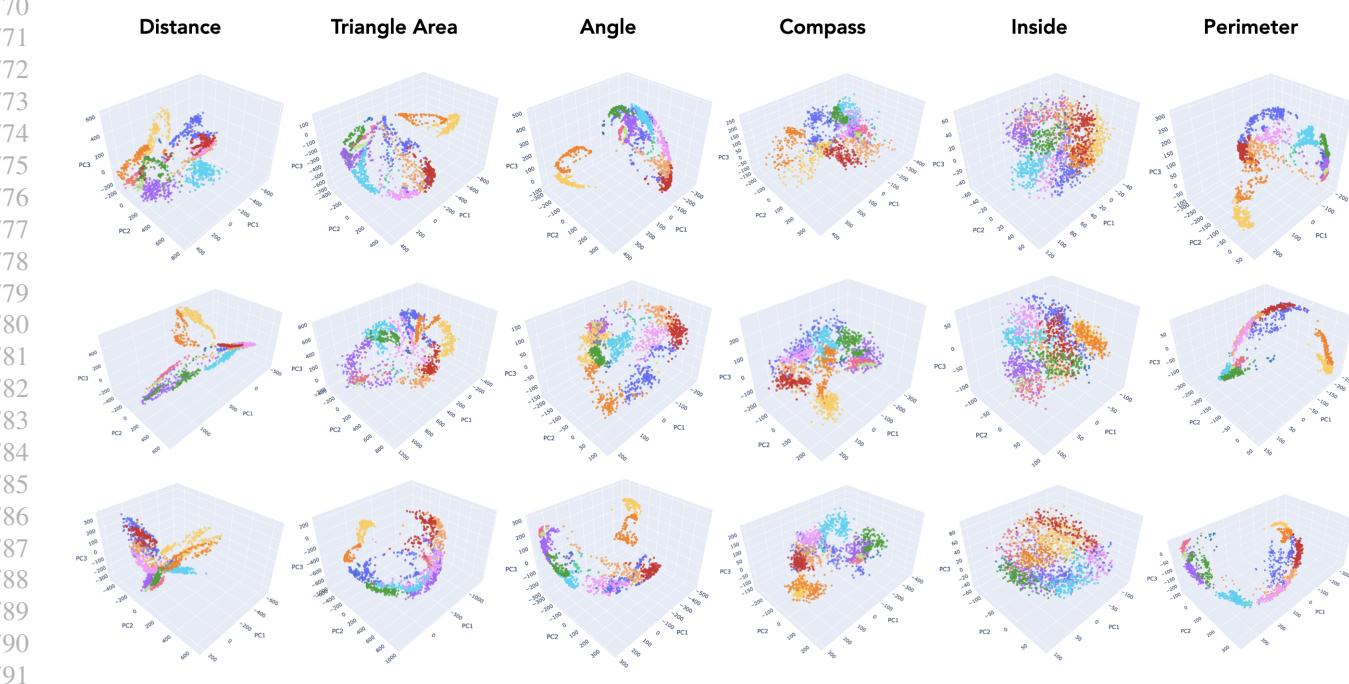


Figure 8. Representation visualizations for single-task models across multiple seeds. Each column shows a different task; each row shows a different random seed. Cities are colored by geographic region.

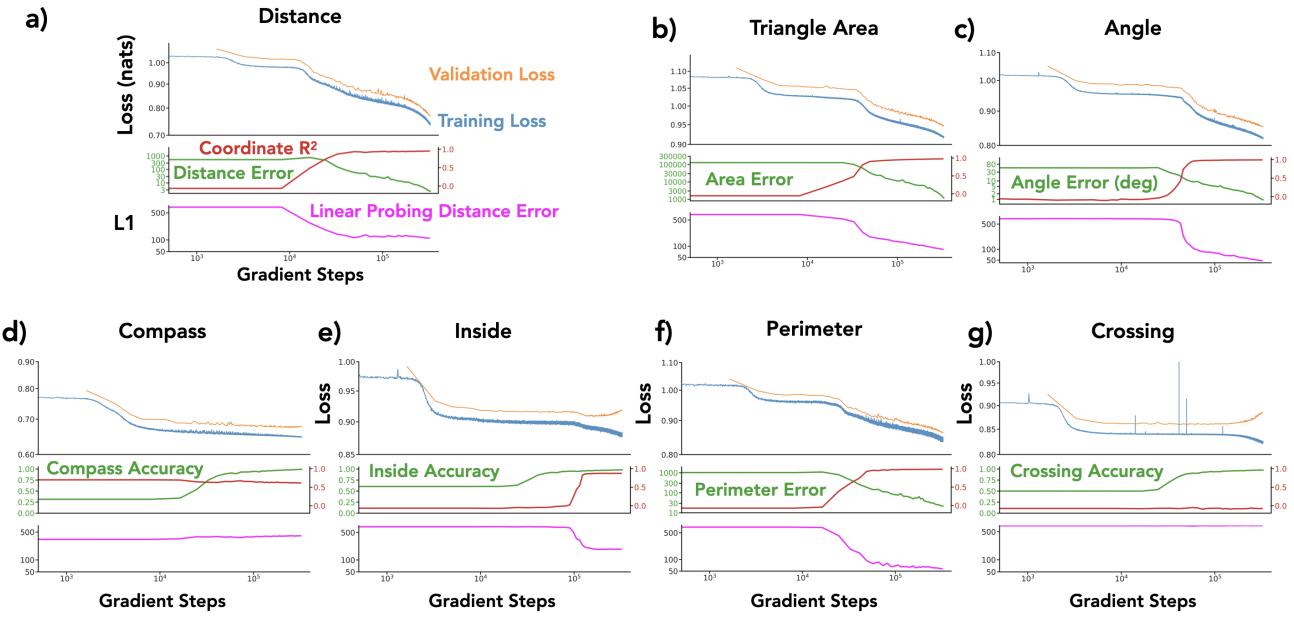


Figure 9. Training dynamics for all single-task models. (a) distance, (b) trianglearea, (c) angle, (d) compass, (e) inside, (f) perimeter, (g) crossing.

CKA vs. Task Count (Per-Seed). Fig. 13 shows the same CKA vs. task count analysis broken down by individual seeds.

Aggregated CKA Trends. Fig. 14 shows CKA vs. task count using all two-task and three-task models.

CKA vs. Generalization (Annotated). Fig. 15 is an annotated version of Fig. 5(b).

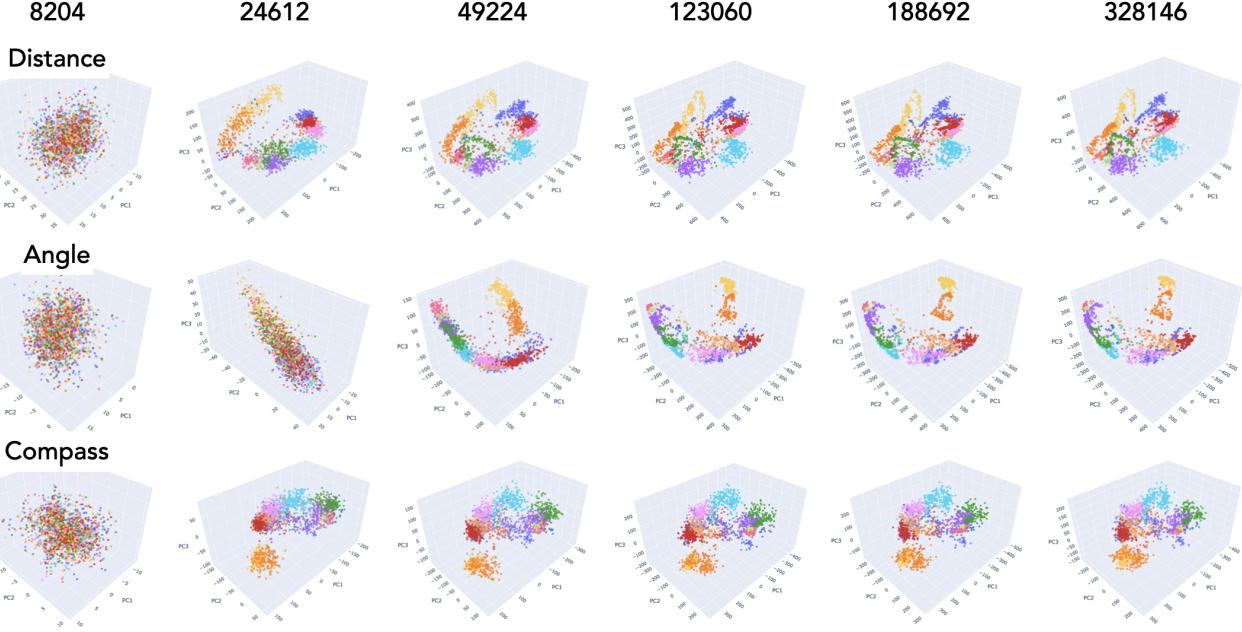


Figure 10. **Representation dynamics during training.** Rows: distance (top), angle (middle), compass (bottom). Columns show PCA projections at gradient steps 8204, 24612, 49224, 123060, 188692, and 328146 (left to right).

C.4. Additional Fine-Tuning Evaluation Results

C.5. Pretraining Variations

Pretraining with Atlantis. Fig. 19 shows the resulting representations when **Atlantis** cities are included from the start of pretraining.

Wider Model. Fig. 20 shows fine-tuning results for a wider model with $2\times$ the hidden dimension.

D. Extended Related Work

See Sec. 2 for main related work.

Interpretability & Internal Representations. Understanding internal representations has roots in neuroscience, informing early neural network development. Beyond the world model discoveries cited in Sec. 2, similar representations emerge during in-context learning. Researchers have also uncovered that models represent meaningful properties of data, concepts, features and abstractions, in interpretable ways.

Fine-tuning. Beyond the works cited in Sec. 2, fine-tuning has been studied extensively across diverse directions: parameter efficiency, zeroth-order optimization, weight composition and representation adaptation.

Dynamics of Representations. Recent work has begun studying how representations evolve during in-context learning or fine-tuning. How internal representations adapt at

inference time is an active area of research.

Loss Plateaus. Our **crossing** task fails to learn in single-task training despite escaping an initial plateau. Such plateaus are notoriously difficult for transformers. Multi-task training has been shown to shorten loss plateaus, similar to why our **crossing** task trains successfully when joined with any other task.

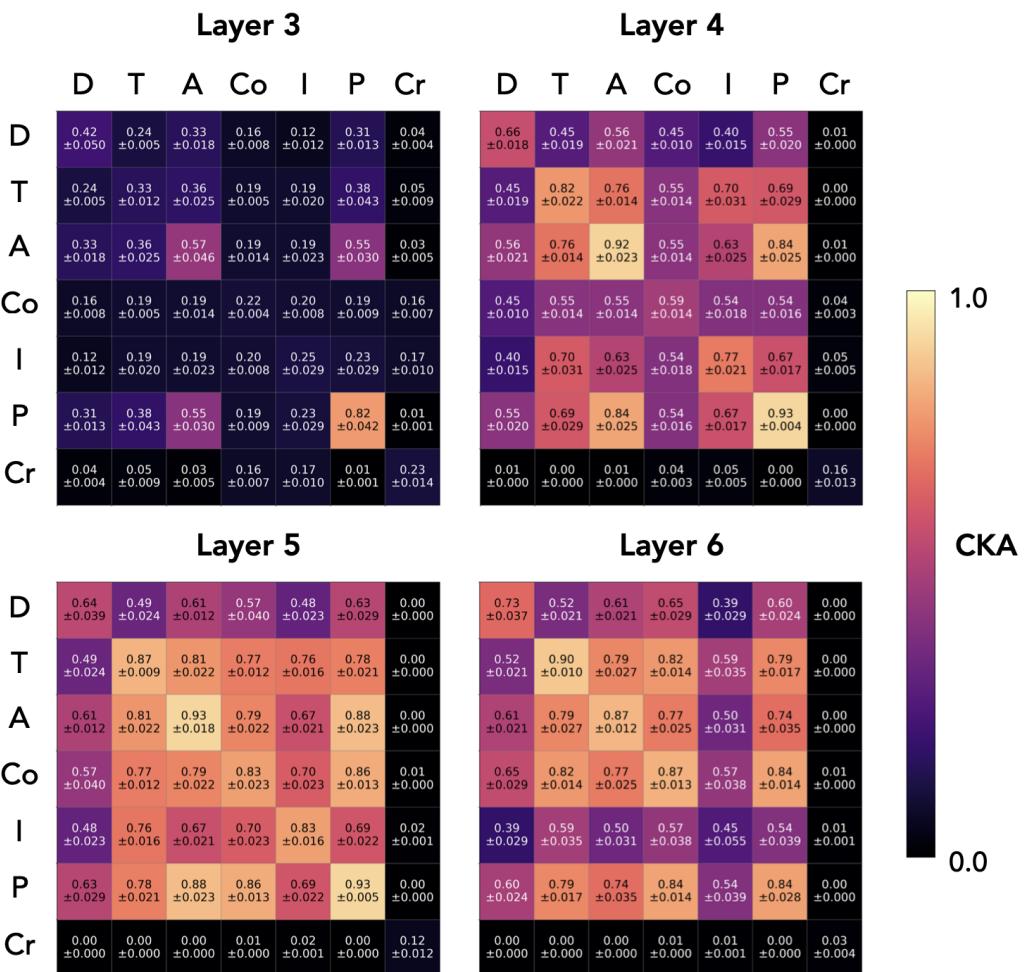
880
881
882883
884
885
886
887
888
889890
891
892
893
894
895
896897
898
899900
901
902
903
904
905
906
907
908
909
910
911912
913
914
915
916
917918
919
920
921
922
923
924
925
926
927
928
929
930931
932
933
934

Figure 11. CKA matrices for single-task models across layers. Each cell shows mean ± SEM across 3 seeds. D=distance, T=triangle area, A=angle, Co=compass, I=inside, P=perimeter, Cr=crossing.

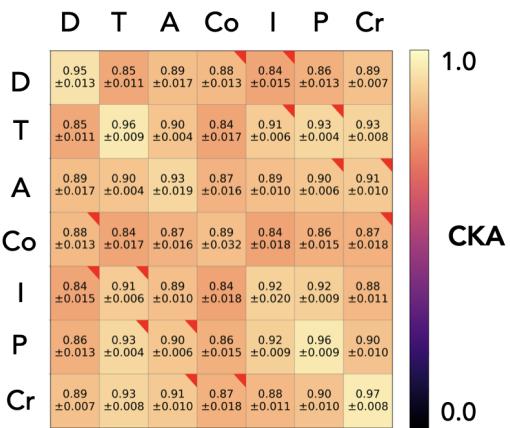


Figure 12. CKA matrix for two-task models at layer 5. Mean ± SEM across 3 seeds. All pairs show high alignment (>0.84).

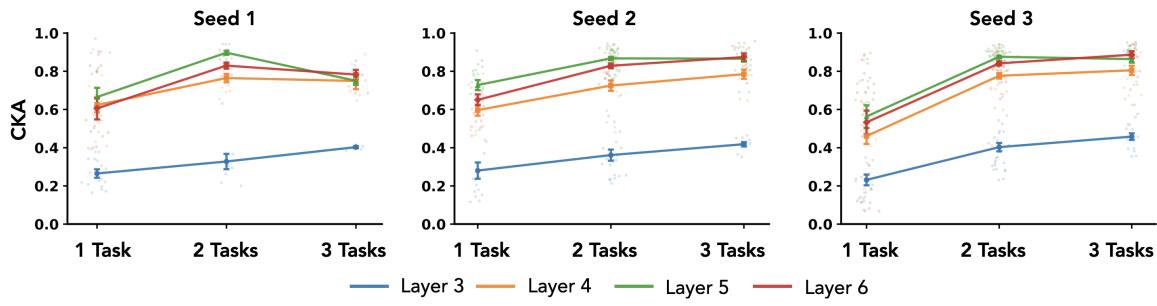


Figure 13. CKA vs. task count for individual seeds. Each panel shows a different seed.

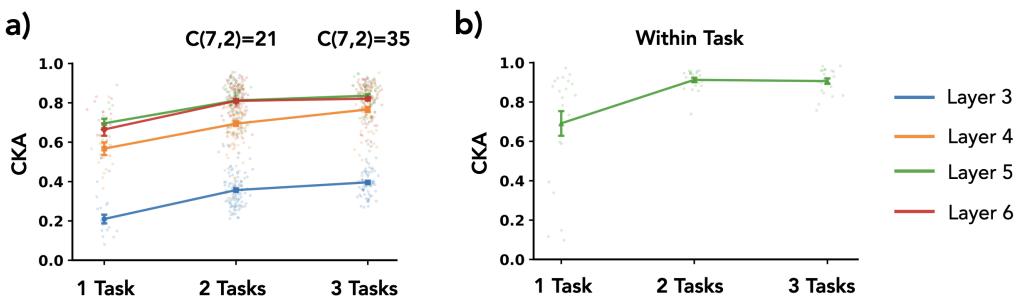


Figure 14. Aggregated CKA analysis. (a) CKA vs. task count for single seed. (b) Within-task CKA increases with task count.

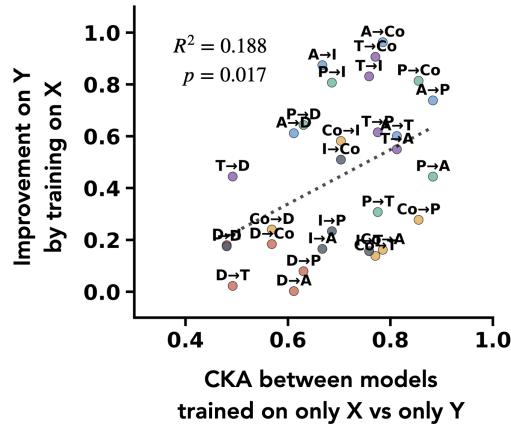


Figure 15. Annotated version of Fig. 5(b). Each point is labeled with its (train→eval) task pair.

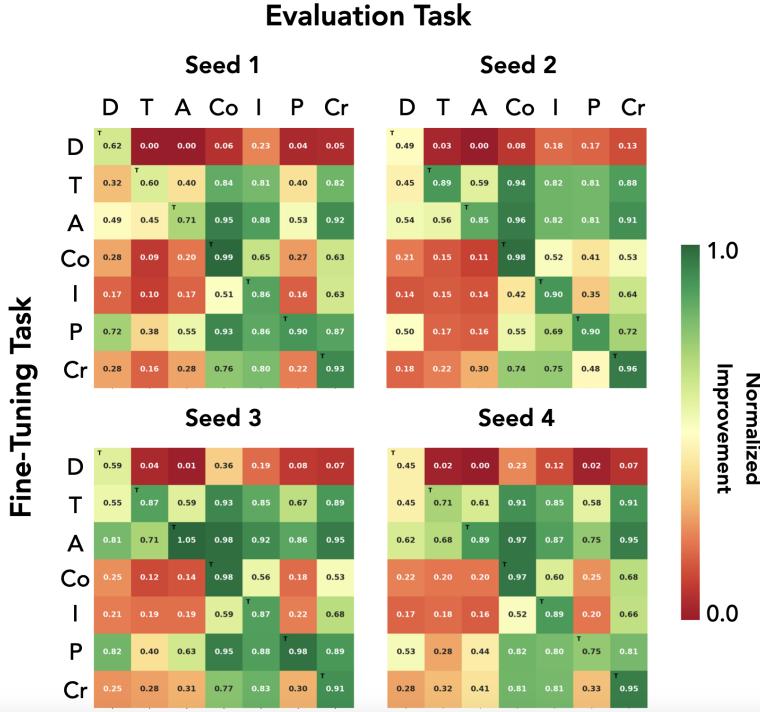


Figure 16. Single-task fine-tuning results for individual seeds. Per-seed version of Fig. 5(a).

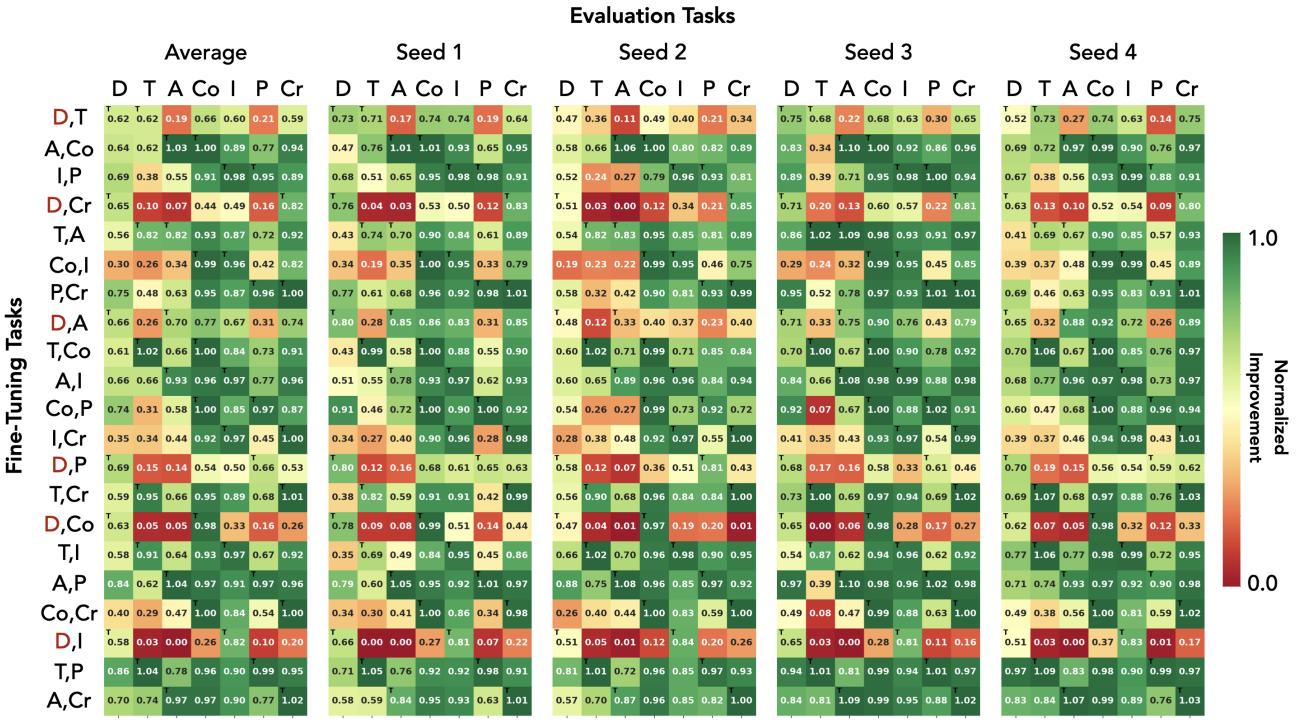


Figure 17. Two-task fine-tuning normalized improvement for all 21 task combinations. Leftmost panel shows average across seeds; remaining panels show individual seeds.

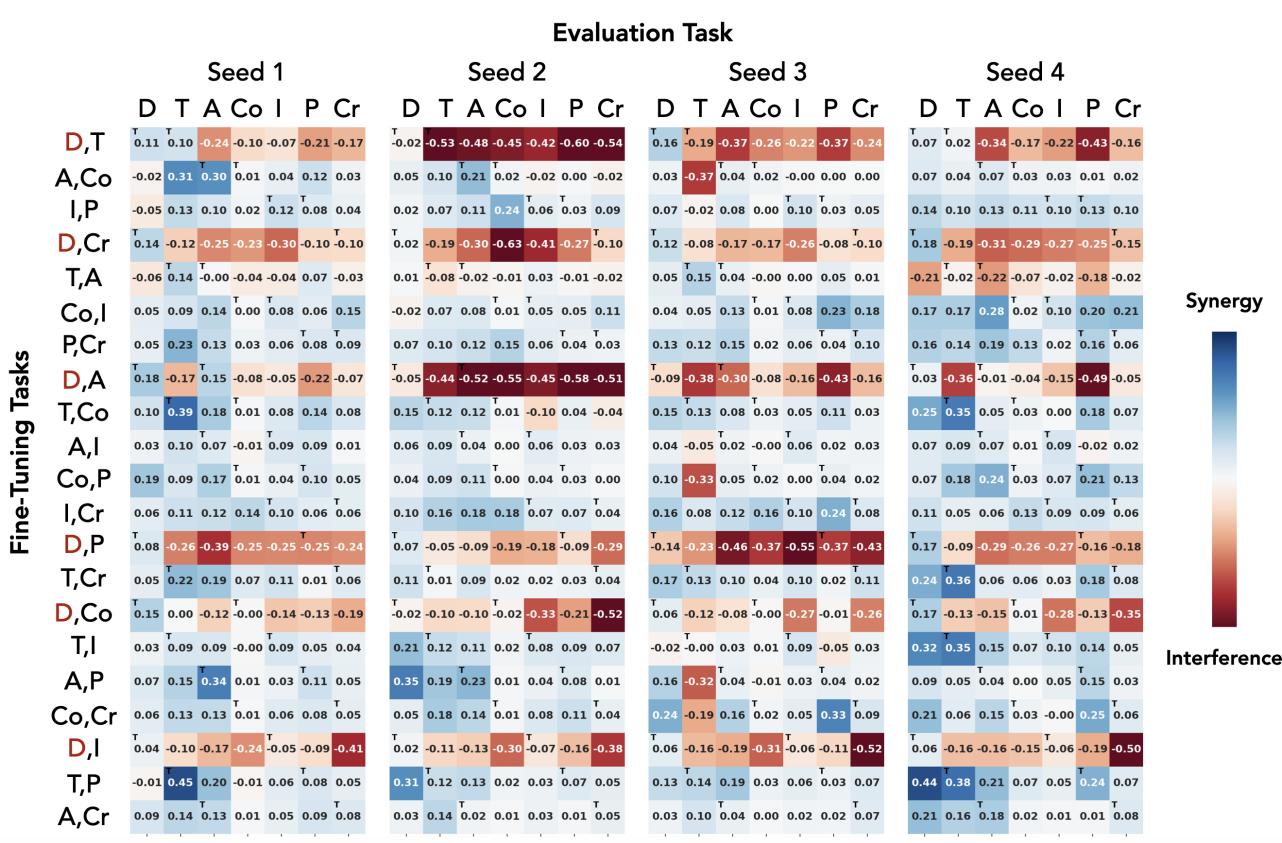
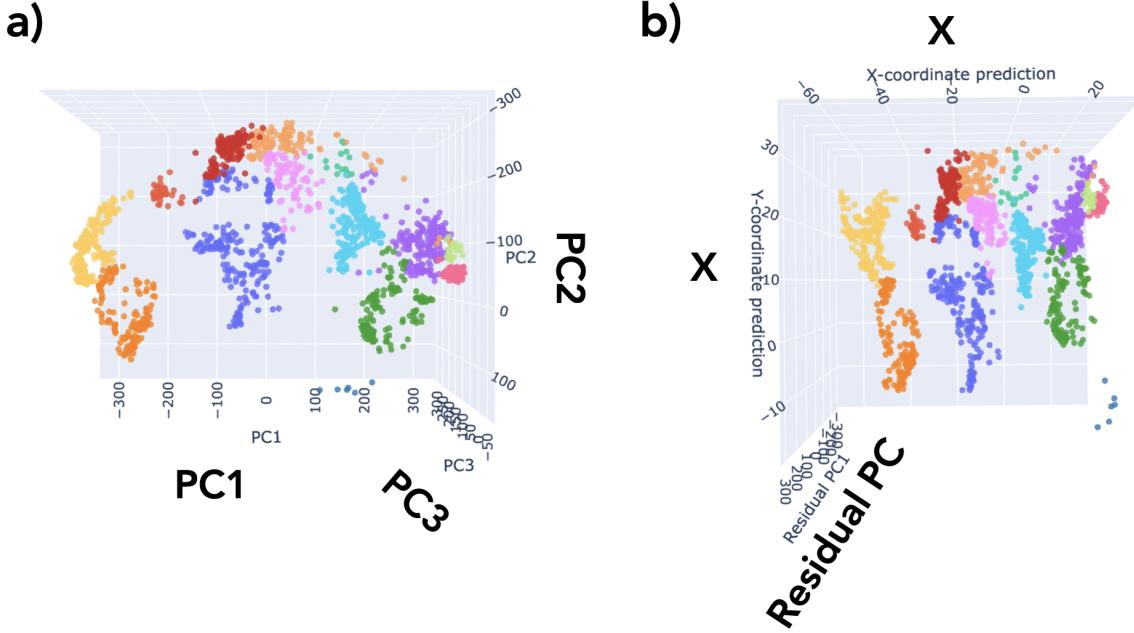


Figure 18. Deviation from best-teacher expectation for all 21 two-task combinations. All 4 seeds shown.

Figure 19. Representations when **Atlantis** is included during pretraining. (a) PCA projection showing **Atlantis** cities integrated with world cities. (b) Linear probe reconstruction confirming geographic accuracy.

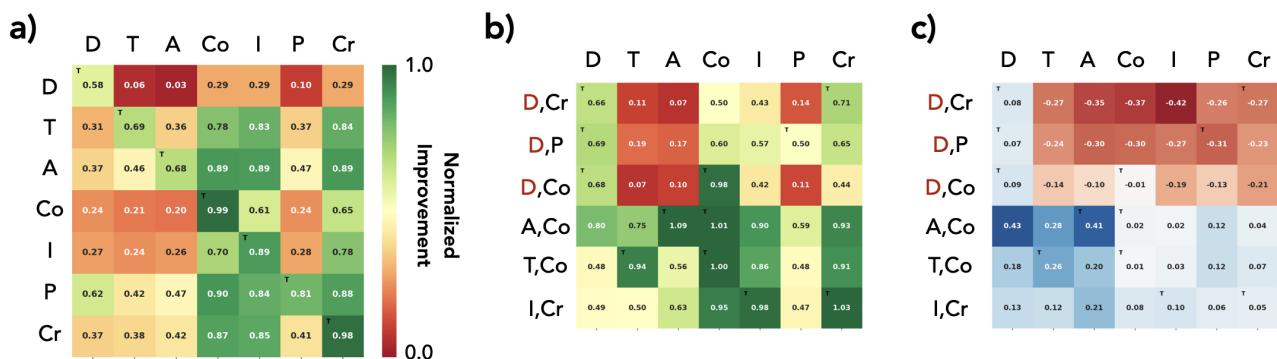


Figure 20. Fine-tuning results for wider model (2× hidden dimension). (a) Single-task fine-tuning normalized improvement. (b) Two-task fine-tuning normalized improvement. (c) Deviation from best-teacher expectation.