

WHEN IS MODEL SOUPING TASTY?

SIMILARITY, TRANSITIVITY, AND ROBUSTNESS

Anonymous authors

Paper under double-blind review

ABSTRACT

Model souping is a post-training technique where the parameters of models are averaged, often leading to improved performance over constituent models without increasing inference cost. However, the specific conditions required for success are not well understood, particularly regarding the trade-off between model diversity and stability. We analyse over 5,000 two-model ResNet-50 soups trained on CIFAR-100, with diversity controlled by branching ingredients from a shared training trajectory at varying epochs. We find that effective souping requires a balance: models must be similar enough to avoid model collapse, but diverse enough to yield improvements. We show that we can predict soupability relatively well with standard similarity metrics. Furthermore, we provide empirical evidence for the hypothesis that souping works by averaging within a low-loss basin by showing that souping is moderately transitive. We also observe that soup gains on corrupted data are strongly correlated with those on in-distribution data. Our findings yield practical advice for machine learning practitioners: if you want a tasty soup, use the right cooking time! Code and experiments are available at: <https://anonymous.4open.science/r/too-salty-478E/>.

1 INTRODUCTION

Wortsman et al. (2022) introduce *souping*: averaging the model parameters produced by different fine-tuning trajectories from a common pre-trained model can yield better generalization than any one ingredient. Souping is formulated as a linear interpolation of model *parameters* θ using *weights* α ; for the case of two models, $\theta_{\text{soup}} = (1 - \alpha)\theta_A + \alpha\theta_B$, creating a new single model. This not only captures the benefit of multiple optimization paths but also introduces an additional adaptation opportunity, as Croce et al. (2023) suggest that dynamically adjusting the interpolation weights enables intermediate behaviours that can better match a range of distributional shifts. In contrast to ensembling which combines the *outputs* of models, souping combines the *parameters*, maintaining the same computational cost of inference with a single forward pass.

Wortsman et al. (2022) hypothesise that souping works because fine-tuned models often lie within the same low-loss basin. A convex combination of their parameters is expected to remain within this basin while reducing the variance introduced by noisy training. They find that when the angle formed by the pre-trained model and ingredients to be souped together is larger, implying greater diversity, there is a greater performance boost from souping. Understanding why averaging maintains a low loss, how diversity contributes to robustness, and when souping is beneficial is therefore essential for studying adaptation more broadly.

Related Work: Souping has been used in a variety of settings. Croce et al. (2023) soup models trained to be robust to different distribution shifts and Ramé et al. (2023) soup ingredients trained on different tasks. Both cases lead to better generalisation. Jang et al. (2025) show that fine-tuned parameters are distributed around a high-accuracy center and use the angle between as few as two fine-tuned weights and their shared branch point to approximate this low-loss center.

Souping is related to the idea of SWA Izmailov et al. (2019). In SWA, the ingredients of the soup come from different steps along the same training trajectory. By contrast, souping averages models from independent training trajectories from a shared initialization.

Souping success requires models to be ‘compatible’ in the sense that averaging their weights has low loss. This notion is closely related to ‘stability to SGD noise’ Frankle et al. (2020), enabled by linear mode connectivity between models trained from the same initialisation under different SGD noise. Such stability is only present after sufficient shared training Iyer et al. (2024). We propose that this defines window during training in which souping is effective and that pre-emptive souping leads to a high loss barrier, or ‘model collapse’.

Our Contributions: Following these works, we seek to better understand souping by conducting a series of experiments addressing the following questions:

How much shared training is required for souping to be effective? We investigate varying the number of shared pre-training epochs before splitting into fine-tuned variants to empirically map the transition from model collapse to effective souping.

Can model similarity predict the effectiveness of souping? Identical models yield no generalization benefit, while overly dissimilar models fail to inhabit a shared low-loss region. We investigate whether standard similarity metrics can predict this balance between model compatibility and diversity.

Is souping transitive? If model A soups with B , and B soups with C , will A soup with C ?

Does souping in-distribution predict souping out-of-distribution? While souping has been shown to help with robustness to distribution shifts, we seek to answer how correlated the soup gains are between in-distribution and out-of-distribution data.

More experiments on the effect of permuting models Ainsworth et al. (2023) prior to souping and how souping affects robustness can be found in Appendix A.1 and Figure A15 respectively.

Soups and Soup Gain: We soup pairs of models using the simple arithmetic mean $\theta_{\text{soup}} = \frac{1}{2}\theta_A + \frac{1}{2}\theta_B$. As shown by Ainsworth et al. (2023) and Iyer et al. (2024), the loss barrier typically has the most extreme behaviour at the midpoint, making it a useful summary statistic of souping performance across weights interpolation, accurately identifying souping failure. Characteristic weighting plots are shown in A16.

Soup gain as the test set accuracy gain of souping two models relative to the *best* parent model.

$$\text{soup gain} = \text{acc}(\theta_{\text{soup}}) - \max\{\text{acc}(\theta_A), \text{acc}(\theta_B)\}$$

where $\text{acc}(\cdot)$ denotes test accuracy. We compare against the maximum accuracy of the ingredients, rather than the mean, as the purpose of a soup should be to improve over its ingredients. Soup gain can alternatively be measured over the loss. We use soup gain as the primary measure of the effectiveness of souping. We also refer to soups with positive and negative soup gain simply as positive and negative soups.

2 EXPERIMENTS

We train a baseline model for image classification on the CIFAR-100 dataset with ResNet-50 He et al. (2016) following Dadalto (2023), saving checkpoints every 10 epochs. From each checkpoint we train 4 new models with different optimizer settings, for details see Appendix Tables 2 and 3. All models are trained to convergence, with the best validation scored model saved for experiments, as shown in Figure A2. A total of 4 variants and 26 branch points were trained, yielding 104 related models and 5,356 binary souping combinations for analysis. Additionally, we train 12 baseline models with and without Stochastic Weight Averaging (SWA) to compare the depth of fine-tuning paths against the breadth of SWA, see Appendix A.4.

2.1 SOUP GAIN AND SHARED EPOCHS

Figure 1 shows the cumulative distribution function (CDF) of soup gains. Many soups have extreme behaviour, with 40% of soups losing over 60% and only 14% with positive gain. There is also a sharp transition to collapse, with only 15% soups between -70% and -20% gain. Grouping by the shared number of training epochs before branching (e.g., ingredients branched at epochs 50 and 100 share 50 epochs), we find the probability of positive soup gain generally increases with shared epochs, reaching around 80% after 260 shared epochs, see Figure 2. At this stage, almost no soups drop accuracy by more than 5%. However, peak gains ($> 0.5\%$) require a balance; if shared training is too extensive ingredients are not diverse enough, performance diminishes. Further analysis on mean gain, loss metrics, and model collapse is provided in Appendices A4, A3 and A5.

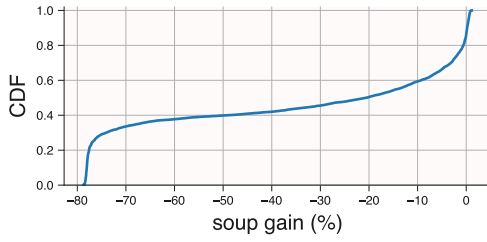


Figure 1: Empirical CDF of soup gain over all soups.

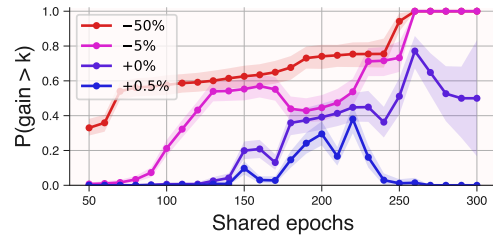


Figure 2: Probability of soup gain being greater than $k\%$ for varying k with 95% CIs.

2.2 PREDICTING SOUPABILITY WITH SIMILARITY

To test if soupability is predictable, we compute a variety of similarity and distance metrics between model pairs. All metrics perform similarly, see Figure A9. Using KL divergence between the outputs of the ingredients as an example, we find a strong negative correlation with soup gain (Spearman -0.86), indicating divergence often leads to model collapse, see Figure 3. However among positive soups, gain has a moderate positive correlation with dissimilarity (Spearman 0.39), see Figure 4. Effective souping requires models to be sufficiently similar, but also rewards variety. Balancing these two effects is key to tasty soups. Additionally, shared epochs correlate strongly with all similarity metrics, for example (Spearman -0.67) with KL divergence, more details shown in Appendix A6.

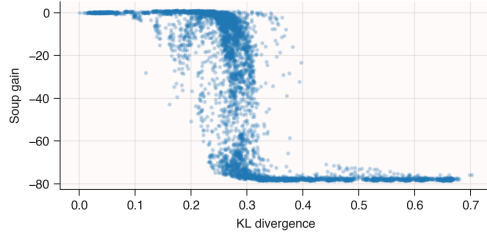


Figure 3: KL vs soup gain (Spearman -0.86).

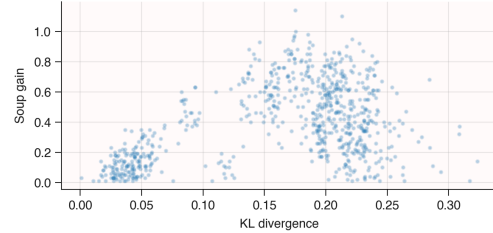


Figure 4: KL vs (+) soup gain (Spearman 0.39).

2.3 IS SOUPING TRANSITIVE?

To test the hypothesis that soupable models reside in the same low-loss basin, we evaluate the transitivity of model triplets (A, B, C) . Figure 5 shows that B and C are highly likely to soup only if A soups with both B and C . We observe a moderate positive correlation (Spearman 0.64) between the gain of (B, C) and the minimum gain of (A, B) and (A, C) , see Figure A11. We attempt to explain when transitivity fails in terms of the branching epochs of the vertices of a given triple, but find no clear pattern. See Appendix Tables 4 and 5.

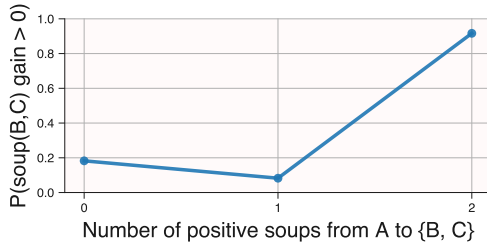


Figure 5: Probability of positive soup gain of B and C vs positive soups with A .

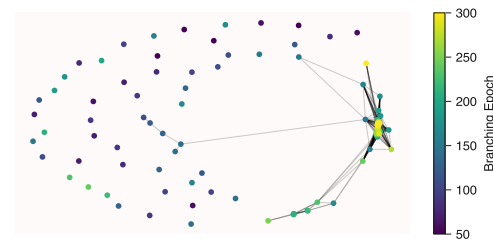


Figure 6: 2D embedding of 104 models using soup-gain distance; edges indicate (+) soups.

We embed our 104 models into 2D using soup gain as a distance metric to search for souping clusters, plotting them in Figure 6, with details found in Section A.6. We find all most successful soups form the edges of a single connected component of models, while there is a single dense cluster of models which are largely branched from later on in the training procedure. There are many counter-examples, but this supports the conclusion that souping is moderately transitive and successful ingredients generally lie within a single loss basin.

2.4 SOUPING FOR ROBUSTNESS TO CORRUPTION

To establish whether souping for in-distribution (ID) performance also increases out-of-distribution (OOD) performance, we compute the soup gain on CIFAR-100C (Hendrycks & Dietterich, 2019) with severity level 3. The soup gains on test and corrupted data correlate strongly (Spearman 0.99), even when restricted to positive soups (Pearson 0.61), see Figures A12, A14. ID performance improvement transfers to unseen target distributions.

We also plot the probability of positive soup gain on corrupted data as a function of shared epochs in Figure A13. The probability of positive gain increases with the number of shared epochs for both clean and corrupted data but the corrupted data always has a lower probability of positive gain.

How does souping compare to SWA? Model accuracy on both clean and corrupted data has been shown to improve with souping. Here we compare the improvements to SWA in Tab. 1. SWA offers a significant boost in robustness to shift, always improves the baseline run, and performs better than our best geometric mean soups. We hypothesize that this is because in this experimental setting, SWA has explored a wider breadth than possible with binary soups and has more closely converged to the low-loss center described by Jang et al. (2025).

Type	P(Gain > 0)	Mean Gain (Acc %)	Mean Corrupted Gain (Acc %)
Soups: 200 Epochs	43.6% \pm 7.7%	0.6% \pm 0.04%	0.71% \pm 0.09%
Soups: 250 Epochs	51.2% \pm 9.9%	0.16% \pm 0.04%	0.05% \pm 0.08%
SWA (12 Runs)	100%	1.6% \pm 0.3%	3.7% \pm 0.4%

Table 1: Comparison of souping robustness to SWA.

3 CONCLUSION

We have tasted more soups to better understand souping. Ingredients must have sufficiently many shared epochs of training in order to be compatible, but not so many that the soup gain is minimal. Various similarity measures between models correlate similarly with soup gain. Similar ingredients are less likely to collapse when souped, but very similar ingredients yield smaller soup gains. Thus, the right balance must be struck for the most effective souping. When souping, we encourage practitioners to test a range of similarities of ingredients to ensure they are finding an optimal soup. Our experiments showing that souping is mostly transitive support the low-loss basin hypothesis. Finally, soup gains on in-distribution data are strongly correlated with those on corrupted data.

Limitations: We investigated ResNet-50 on CIFAR-100 to enable a comprehensive combinatorial analysis of over 5,000 binary soups that would be computationally prohibitive with large foundation models. Iyer et al. (2024) demonstrated that the linear mode connectivity is a fundamental optimization property, suggesting our findings would transfer to larger-scale fine-tuning settings. Additionally, we only consider pairwise souping using the arithmetic mean at the midpoint for similar reasons. Other methods of souping, such as learned soups, may yield different results.

Future Work: Future empirical work could conduct similar experiments in different settings, such as a variety of model architectures and datasets. Theory could be developed for souping in simpler settings like an overparameterized linear model or a shallow network. Theory could also be created to help predict the change in loss we expect from souping associated with noise reduction. Validating the performance of SWA against model stock estimates of the low-loss center could confirm if model souping remains a strong contender for test-time robustness under distribution shift.

REFERENCES

- Samuel K. Ainsworth, Jonathan Hayase, and Siddhartha Srinivasa. Git re-basin: Merging models modulo permutation symmetries, 2023. URL <https://arxiv.org/abs/2209.04836>.
- Francesco Croce, Sylvestre-Alvise Rebuffi, Evan Shelhamer, and Sven Gowal. Seasoning model soups for robustness to adversarial and natural distribution shifts, 2023. URL <https://arxiv.org/abs/2302.10164>.
- Eduardo Dadalto. Resnet-50 model trained on cifar-100. https://huggingface.co/edadaltocg/resnet50_cifar100, 2023. Hugging Face Model Repository.
- Jonathan Frankle, Gintare Karolina Dziugaite, Daniel M. Roy, and Michael Carbin. Linear mode connectivity and the lottery ticket hypothesis, 2020. URL <https://arxiv.org/abs/1912.05671>.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep Residual Learning for Image Recognition. In *Proceedings of 2016 IEEE Conference on Computer Vision and Pattern Recognition, CVPR ’16*, pp. 770–778, Las Vegas, NV, USA, June 2016. IEEE. doi: 10.1109/CVPR.2016.90. URL <http://ieeexplore.ieee.org/document/7780459>.
- Dan Hendrycks and Thomas Dietterich. Benchmarking neural network robustness to common corruptions and perturbations, 2019. URL <https://arxiv.org/abs/1903.12261>.
- Gaurav Iyer, Gintare Karolina Dziugaite, and David Rolnick. Linear weight interpolation leads to transient performance gains. *Transactions on Machine Learning Research*, 2024. URL <https://openreview.net/forum?id=XGAdBXlFcj>. Presented at HiLD (ICML 2024 Workshop).
- Pavel Izmailov, Dmitrii Podoprikin, Timur Garipov, Dmitry Vetrov, and Andrew Gordon Wilson. Averaging weights leads to wider optima and better generalization, 2019. URL <https://arxiv.org/abs/1803.05407>.
- Dong-Hwan Jang, Sangdoo Yun, and Dongyoon Han. Model stock: All we need is just a few fine-tuned models, 2025. URL <https://arxiv.org/abs/2403.19522>.
- Alexandre Ramé, Kartik Ahuja, Jianyu Zhang, Matthieu Cord, Léon Bottou, and David Lopez-Paz. Model ratatouille: Recycling diverse models for out-of-distribution generalization, 2023. URL <https://arxiv.org/abs/2212.10445>.
- Mitchell Wortsman, Gabriel Ilharco, Samir Yitzhak Gadre, Rebecca Roelofs, Raphael Gontijo-Lopes, Ari S. Morcos, Hongseok Namkoong, Ali Farhadi, Yair Carmon, Simon Kornblith, and Ludwig Schmidt. Model soups: averaging weights of multiple fine-tuned models improves accuracy without increasing inference time, 2022. URL <https://arxiv.org/abs/2203.05482>.

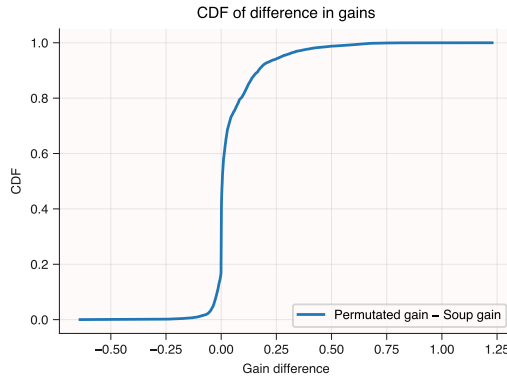


Figure A1: Cumulative distribution function (CDF) of the difference in soup gain before and after permutation alignment using `rebasin`. This ignores the 7% of soups with a loss higher than 5 after permutation as these make the plot difficult to interpret. The remaining mean and median difference is approximately zero, indicating that permutation alignment does not have a significant effect on the effectiveness of souping in our experiments. While some soups benefit from permutation alignment, others are negatively affected, leading to an overall negligible impact while there is a risk of severe degradation.

A APPENDIX

A.1 PERMUTATION ALIGNMENT FOR SOUPING

Following the work from Ainsworth et al. (2023), we investigate whether permuting the neurons of models prior to souping increases the effectiveness of souping. We use the `rebasin`¹ package to align pairs of models before souping. This package uses the ‘matching weights’ method which permutes the neurones by inspecting only the weights. This contrasts with ‘activation matching’ which requires forward passes through the network, and ‘straight through estimators’ which are even more computationally expensive. The authors find that matching weights performs similarly to activation matching while being computationally cheaper. Therefore, we only consider the matching weights method.

We align all 5,346 pairs of models using `rebasin` and compute the soup gain after alignment. Prior to permutation, 14.25% of soups were positive, while after permutation, 14.32% of soups were positive. However, 7% of soups obtained a loss higher than 5, which is worse than the loss of any previous soup. We plot the cumulative distribution function (CDF) of the difference in soup gain before and after permutation in Figure A1. This plot shows that while permuting can sometimes help, it does not do so consistently. Further, the median difference in soup gain is approximately zero, indicating that permuting does not have a significant effect on the effectiveness of souping in our experiments. There also remains a significant risk of severe degradation. Thus, we conclude that ‘matching weights’ permutation does not make a noticeable difference to soupability in our setting.

¹<https://pypi.org/project/rebasin/>

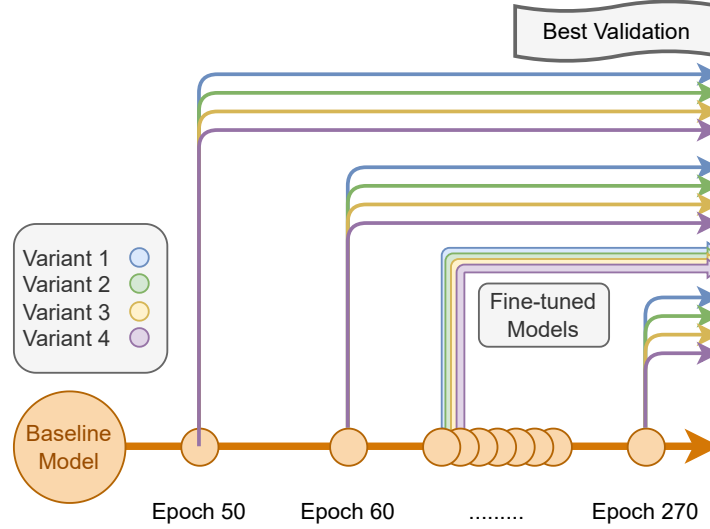


Figure A2: Branching of fine-tuned models from baseline checkpoints. A single baseline model is trained, with checkpoints saved every 10 epochs. From each checkpoint, 4 variants are trained with different optimizer hyper-parameter perturbations.

Model	Learning Rate Scale	Momentum Scale	Weight Decay Scale
Model 1	0.7073	1.1009	0.8284
Model 2	1.2244	0.9247	1.1150
Model 3	0.5112	1.0695	1.1099
Model 4	0.5373	0.8594	0.9078

Table 2: Optimizer perturbation scales applied during finetuning from the baseline ResNet-50 checkpoint on CIFAR-100. Each model scales the original SGD hyperparameters multiplicatively.

A.2 TRAINING DETAILS FOR CIFAR-100 WITH RESNET-50

Component	Hyperparameter	Value
Dataset	Dataset	CIFAR-100
	# Classes	100
	Data augmentation	Mirroring and Padded Offset
	Validation split	5% of training set
	Split seed	42
Model	Architecture	ResNet-50
	Pretrained	No (from scratch)
Optimization	Optimizer	SGD (Nesterov)
	Initial learning rate	0.1
	Momentum	0.9
	Weight decay	5×10^{-4}
Learning rate schedule	Scheduler	CosineAnnealingLR
	T_{\max}	280 epochs
	η_{\min}	0
Training	Epochs	300
	Batch size	128
	Mixed precision	No

Table 3: Training hyperparameters for CIFAR-100 with ResNet-50. A randomization seed of 42 is used unless otherwise specified.

A.3 FURTHER DETAILS ON EXPERIMENTS

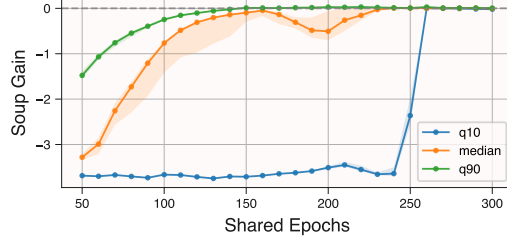


Figure A3: Quantiles of soup gain vs shared epochs with soup gain measured as the reduction in loss vs the best parent. Similar to the conclusions from 2, we find model collapse almost never occurs past shared epochs 250.

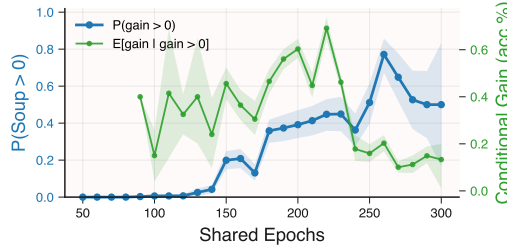


Figure A4: Probability of positive soup gain and conditional expected gain vs shared epochs. We find that the expected soup gain noisy, but is maximised in an ideal window of shared epochs. Past this point, and models are too similar, leading to minimal gains. Before this point, models are incompatible and rarely yield positive gain.

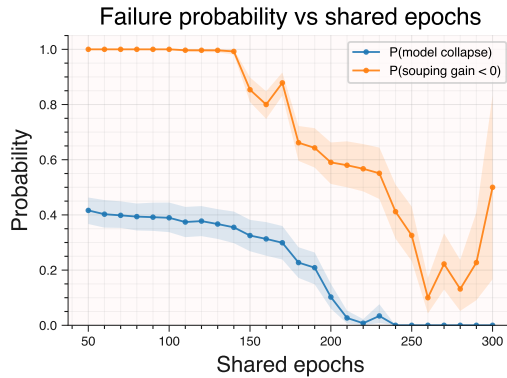


Figure A5: We plot both the probability of model collapse and the probability that soup gain is positive over varying shared epochs, with 95% bootstrapped confidence intervals. Model collapse is defined as the soup attaining less than 5% test accuracy. Both model collapse and negative soup gain decrease with shared epochs. In particular, model collapse is very rare past epoch 200 and never occurs past epoch 250.

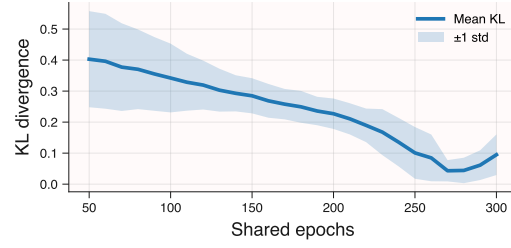


Figure A6: Shared epochs vs KL divergence. We plot the mean and standard deviation over each value of shared epochs. The full data has a Spearman correlation of -0.67 . We conclude that we can noisily recover the number of shared epochs by measuring the KL divergence, and that model similarity can be controlled imprecisely by varying the number of shared epochs.

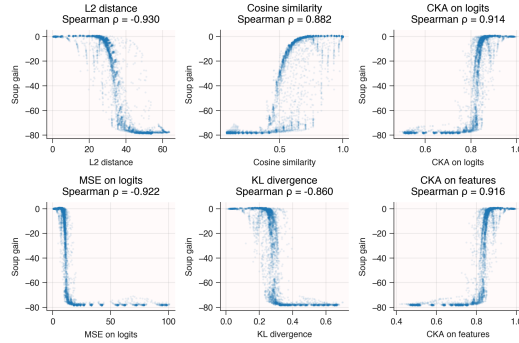


Figure A7: Soup gain vs various similarity and distance metrics between pairs of models. Each subplot shows the soup gain against one metric, with the Spearman correlation. All metrics perform similarly. We conclude that the more similar models are, the more likely souping is to not cause model collapse.

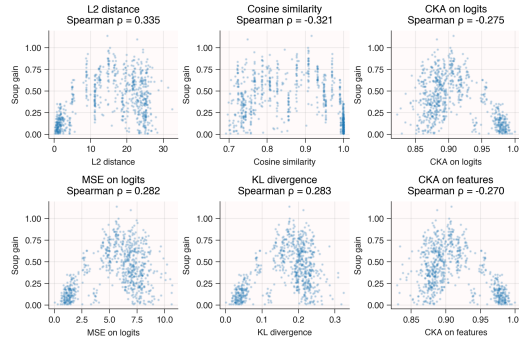


Figure A8: Soup gain vs various similarity and distance metrics between pairs of models, for the subset of soups with positive soup gain. Each subplot shows the soup gain against one metric, with the Spearman correlation. Each metric correlates similarly with soup gain. We see that when models are very similar, soup gain is small, while more dissimilar models can yield larger soup gains.

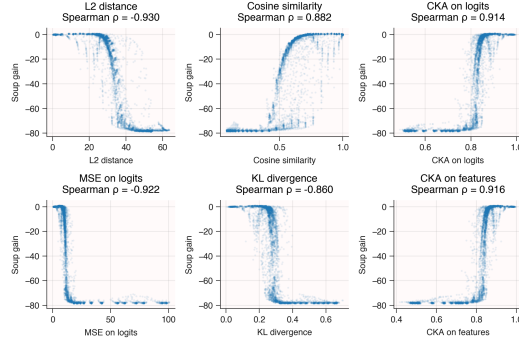


Figure A9: Soup gain vs various similarity and distance metrics between pairs of models. Each subplot shows the soup gain against one metric, with the Spearman correlation. All metrics perform similarly. We conclude that the more similar models are, the more likely souping is to not cause model collapse.

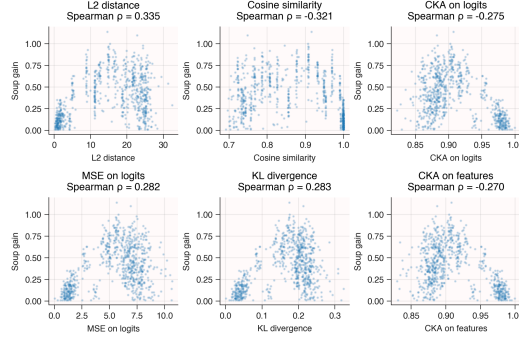


Figure A10: Soup gain vs various similarity and distance metrics between pairs of models, for the subset of soups with positive soup gain. Each subplot shows the soup gain against one metric, with the Spearman correlation. Each metric correlates similarly with soup gain. We see that when models are very similar, soup gain is small, while more dissimilar models can yield larger soup gains.

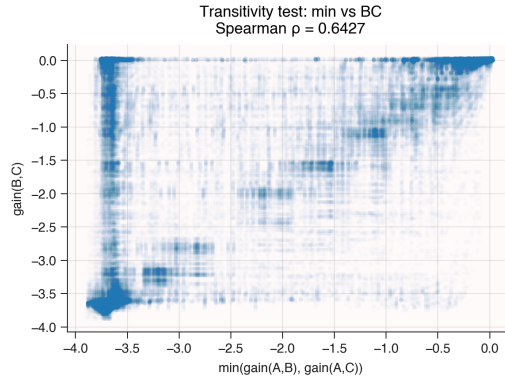


Figure A11: Scatterplot of the soup gain of models B and C against the minimum soup gain of models A with B and C . Each point represents a triplet of models (A, B, C) . We observe a positive Spearman correlation of 0.64, suggesting that souping is fuzzily transitive. The correlation with the mean soup gain is lower, at 0.49.

Table 4: Transitivity failure vs shared epochs over all (A, B, C) triples. Here n denotes the number of triples in the bin for which two soups are positive, and p_{fail} is the proportion of those triples for which the remaining soup is negative. Failure rates are lowest for very small and very large shared prefixes, with a peak at intermediate shared epochs, but no monotonic trend is observed.

Shared epochs (ABC)	n	p_{fail}
150–180	2294	0.071
180–200	1887	0.075
200–220	1657	0.119
220–240	1244	0.104
240–300	1369	0.055

Table 5: Transitivity failure vs branching-epoch difference $|B - C|$. Here n and p_{fail} are defined as in Table 4. We restrict to triples with shared epochs in $[200, 220]$, fix the lowest-epoch model, and vary the branching-epoch difference of the remaining two models. We observe no clear dependence of transitivity failure on epoch difference within this regime. Similar trends were observed for the other bins which are omitted for brevity.

Epoch diff (BC)	n	p_{fail}
0–10	792	0.124
10–20	497	0.123
20–30	419	0.124
30–50	578	0.080
50–100	369	0.100

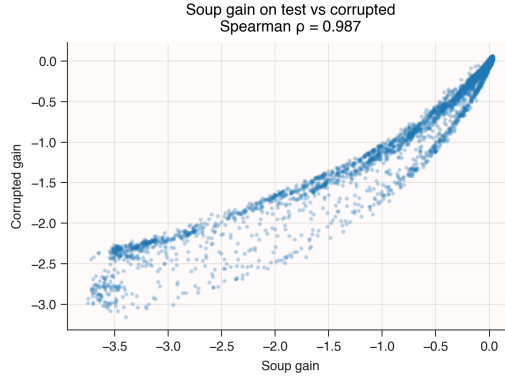


Figure A12: Scatterplot of the soup gain on test vs corrupted data. There is a clear sub-linear trend with strong correlation. Such a close relationship is sensitive to the nature of the distribution shift. This plot mostly shows that when model collapse occurs on the original test set, it also occurs on the corrupted data.

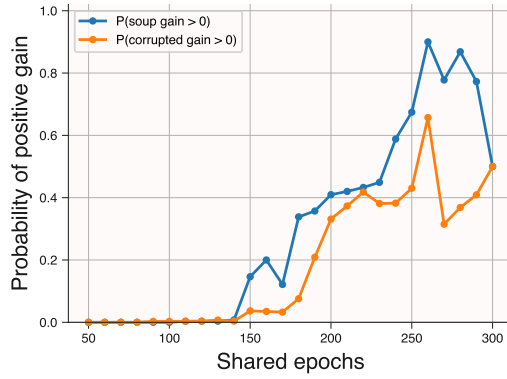


Figure A13: Probability of positive gain for soups as a function of the number of shared epochs. We see that the probability of positive gain increases with the number of shared epochs, for both clean and corrupted data. However, the corrupted data consistently has a slightly lower probability of positive gain. Thus, while souping also helps on corrupted data, it is slightly less effective than on clean data.

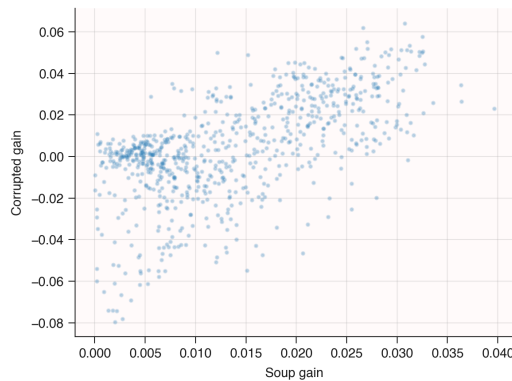


Figure A14: Plot of soup gain on test vs corrupted data for only models with positive soup gain on the test set. Spearman correlation of 0.61.

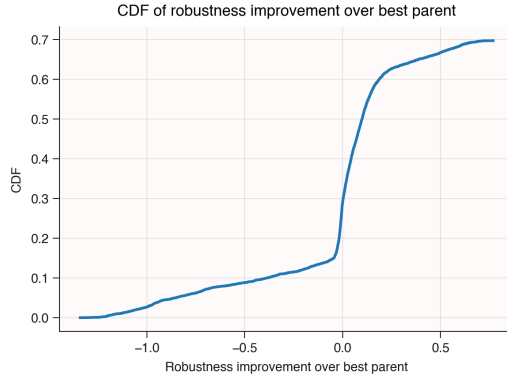


Figure A15: CDF of the difference in *robustness gap* before and after souping. The robustness gap is defined as the difference in loss between the test set and the corrupted set. The robustness gap before souping is taken as the minimum robustness gap of the parents. We see a fairly symmetric distribution. It has mean -0.04 and median 0.02. We therefore conclude that souping does not systematically improve the robustness gap.

A.4 SWA EXPERIMENT RUNS

Baseline training was performed using the same hyper-parameters as in 3 and using the top validation model. We compare to SWA starting at epoch 220 from the same base model. The initialization and batch ordering seeds with full results are shown below in 6.

Seed	Clean Data						Corrupted Data					
	Accuracy (\uparrow)			Loss (\downarrow)			Accuracy (\uparrow)			Loss (\downarrow)		
	Base	SWA	Δ	Base	SWA	Δ	Base	SWA	Δ	Base	SWA	Δ
42	78.39	79.84	+1.45	1.02	0.90	-0.12	50.88	54.05	+3.17	2.37	2.63	+0.26
43	78.04	79.81	+1.77	1.03	0.87	-0.16	50.82	54.47	+3.65	2.35	2.51	+0.16
44	78.87	80.64	+1.77	1.00	0.86	-0.14	51.31	55.26	+3.95	2.36	2.57	+0.21
45	78.57	80.18	+1.61	1.00	0.86	-0.14	50.76	54.59	+3.83	2.39	2.60	+0.21
46	78.04	79.42	+1.38	1.02	0.89	-0.13	51.06	54.46	+3.40	2.35	2.57	+0.22
47	78.47	80.38	+1.91	1.01	0.88	-0.13	51.51	55.28	+3.77	2.35	2.56	+0.21
48	78.21	79.85	+1.64	1.02	0.90	-0.12	50.37	54.29	+3.92	2.40	2.62	+0.22
49	78.57	80.27	+1.70	1.01	0.88	-0.13	51.62	56.05	+4.43	2.33	2.44	+0.11
50	78.97	79.88	+0.91	1.00	0.89	-0.11	51.59	55.00	+3.41	2.34	2.56	+0.22
51	78.09	79.87	+1.78	1.03	0.89	-0.14	51.52	55.25	+3.73	2.32	2.52	+0.20
52	78.48	79.86	+1.38	1.02	0.89	-0.13	51.05	54.22	+3.17	2.37	2.65	+0.28
53	78.55	80.26	+1.71	1.02	0.88	-0.14	51.77	55.42	+3.65	2.31	2.50	+0.19
Mean	78.44	80.02	+1.58	1.01	0.88	-0.13	51.19	54.86	+3.67	2.35	2.56	+0.21
Std	0.30	0.33	0.27	0.01	0.01	0.01	0.43	0.60	0.36	0.03	0.06	0.04

Table 6: Comparison of baseline and SWA models branching from Epoch 220.

A.5 BINARY MODEL SOUP ACCURACY CURVES

While our experiments make a statistical analysis on the characteristics of even-weighted binary soups, we verify the validity of our simplifying assumptions with characteristic accuracy curves across 16 binary pairs to support our claim that even-weighted soups are a good measure for performance. For soups exhibiting model collapse, the measure is much more accurate, while for positive gain soups the metric is indicative of soupability but not as accurate for actual souping gain.

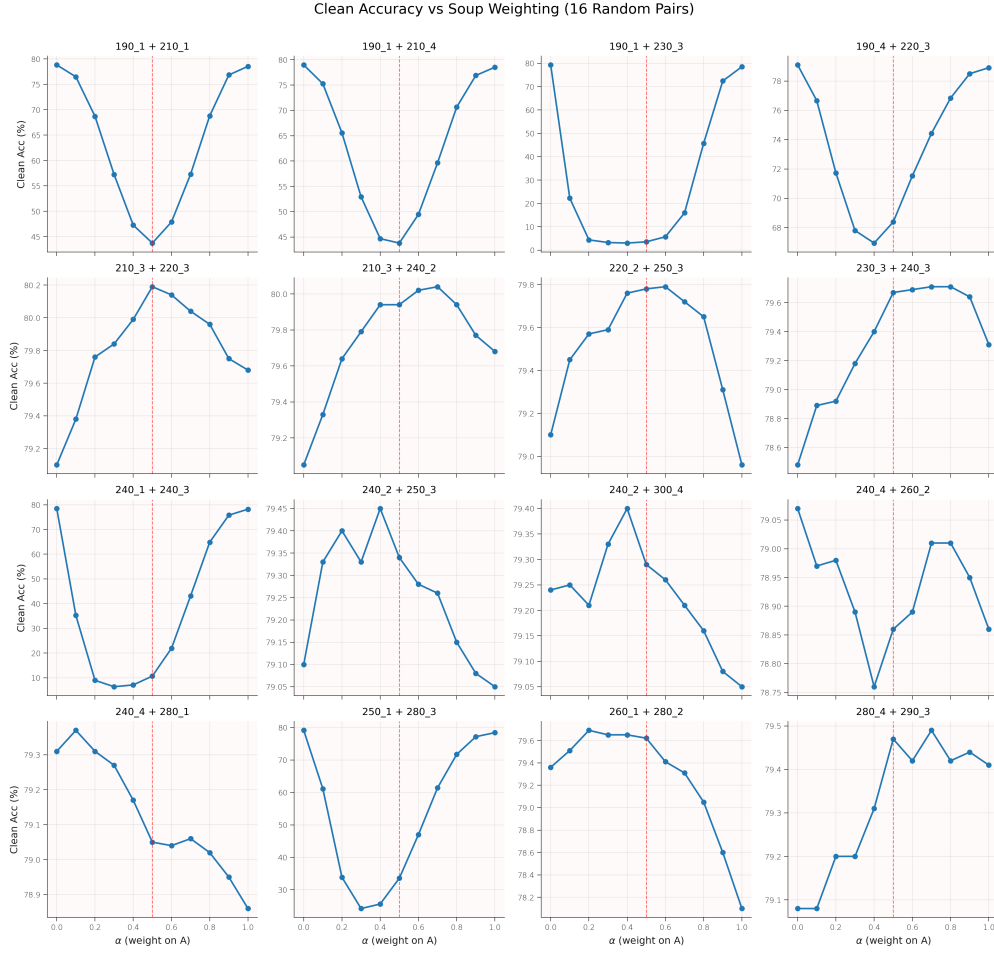


Figure A16: Souping characteristics across a range of parameter interpolation values for 16 models.

A.6 DETAILS FOR MODEL EMBEDDING CREATION

Given our evidence for transitivity, we consider the broader landscape of all 104 of our originally trained models. Do the soups all exist in separate clusters of models in separate loss basins? We define a distance metric defined as

$$d_{AB} = -\text{sign}(\text{soup gain}) - 0.1 * \text{soup gain}$$

where d_{AB} is the distance between models A and B . Intuitively, this metric puts models close together that soup together positively, taking into account the magnitude of soup gain. We then cast this down into a 2-dimensional embedding using Multidimensional Scaling. We also color by branching epoch and mark edges that represent successful soups. The resulting plot is shown in Figure 6.