Infrastructure and Progress Towards the First Community-Built and Continually-Improved Model Colin Raffel

# Collaborative, Communal, Continual {CCC} ML

- Where are we in terms of making CCCML possible?
- We can build a useful VCS for CCCML now.
- We're building a git-based VCS for CCCML.
- We're building a CCC model to test it out.



































from <u>https://dvc.org/</u>

How can we enable collaborative and continual development of machine learning models?

We need to be able to cheaply communicate **patches** and **merge** updates from different contributors. How can we enable collaborative and continual development of machine learning models?

# We need to be able to cheaply

communicate patches and merge updates

from different contributors.



From *Training Neural Networks with Fixed Sparse Masks* by Sung et al.



From Few-Shot Parameter-Efficient Fine-Tuning is Better and Cheaper than In-Context Learning by Liu et al.

$$\mathbf{h}_{2}^{\ell} = \operatorname{Dropout}\left(\mathbf{W}_{m_{1}}^{\ell} \cdot \mathbf{h}_{1}^{\ell} + \mathbf{b}_{m_{1}}^{\ell}\right) \quad (1)$$

$$\mathbf{h}_{3}^{\ell} = \mathbf{g}_{LN_{1}}^{\ell} \odot \frac{(\mathbf{h}_{2}^{\ell} + \mathbf{x}) - \mu}{\sigma} + \mathbf{b}_{LN_{1}}^{\ell} \quad (2)$$

$$\mathbf{h}_{4}^{\ell} = \operatorname{GELU}\left(\mathbf{W}_{m_{2}}^{\ell} \cdot \mathbf{h}_{3}^{\ell} + \mathbf{b}_{m_{2}}^{\ell}\right) \quad (3)$$

$$\mathbf{h}_{5}^{\ell} = \operatorname{Dropout}\left(\mathbf{W}_{m_{3}}^{\ell} \cdot \mathbf{h}_{4}^{\ell} + \mathbf{b}_{m_{3}}^{\ell}\right) \quad (4)$$

$$\operatorname{Dut}^{\ell} = \mathbf{g}_{LN_{2}}^{\ell} \odot \frac{(\mathbf{h}_{5}^{\ell} + \mathbf{h}_{3}^{\ell}) - \mu}{\sigma} + \mathbf{b}_{LN_{2}}^{\ell} \quad (5)$$

From BitFit: Simple Parameter-efficient Fine-tuning for Transformer-based Masked Language-models by Ben Zaken et al.



From LoRA: Low-Rank Adaptation of Large Language Models by Hu et al.



From Measuring the Intrinsic Dimension of Objective Landscapes by Li et al.

**Algorithm 1:** Sparse Ternary Compression (STC) 1 input: flattened tensor  $T \in \mathbb{R}^n$ , sparsity p 2 output: sparse ternary tensor  $T^* \in \{-\mu, 0, \mu\}^n$  $\mathbf{3} \bullet k \leftarrow \max(np, 1)$ 4 •  $v \leftarrow \operatorname{top}_k(|T|)$ 5 • mask  $\leftarrow (|T| \ge v) \in \{0, 1\}^n$ 6 •  $T^{masked} \leftarrow mask \odot T$ 7 •  $\mu \leftarrow \frac{1}{k} \sum_{i=1}^{n} |T_i^{masked}|$ 8 return  $T^* \leftarrow \mu \times \operatorname{sign}(T^{masked})$ 

From Robust and Communication-Efficient Federated Learning from Non-IID Data by Sattler et al.

How can we enable collaborative and continual development of machine learning models?

We need to be able to cheaply communicate **patches** and **merge** updates from different contributors.



From Merging Models with Fisher-Weighted Averaging by Matena and Raffel





From Robust and Communication-Efficient Federated Learning from Non-IID Data by Sattler et al.



From Branch-Train-Merge: Embarrassingly Parallel Training of Expert Language Models by Li et al.



From *Patching open-vocabulary models by interpolating weights* by Ilharco et al.



From *Fusing finetuned models for better pretraining* by Choshen et al.



From Merging Models with Fisher-Weighted Averaging by Matena and Raffel

![](_page_33_Picture_0.jpeg)

From Git Re-Basin: Merging Models modulo Permutation Symmetries by Ainsworth et al.

![](_page_34_Figure_0.jpeg)

How can we enable collaborative and continual development of machine learning models?

We need a system to **track changes** and manage contributions to a model.

- **Track changes to a checkpoint.** There are many ways to modify a checkpoint, e.g. updating all parameters, updating a subset of parameters, adding/removing parameters, etc. The system should support all of these operations.

- **Track changes to a checkpoint.** There are many ways to modify a checkpoint, e.g. updating all parameters, updating a subset of parameters, adding/removing parameters, etc. The system should support all of these operations.
- **An on-disk format for changes.** This format should only store the information required to reconstruct the change and should support all update types.

- **Track changes to a checkpoint.** There are many ways to modify a checkpoint, e.g. updating all parameters, updating a subset of parameters, adding/removing parameters, etc. The system should support all of these operations.
- **An on-disk format for changes.** This format should only store the information required to reconstruct the change and should support all update types.
- **Restoring the previous state of a checkpoint.** The system should be able to rapidly restore the state of the model at a certain point of history.

- **Track changes to a checkpoint.** There are many ways to modify a checkpoint, e.g. updating all parameters, updating a subset of parameters, adding/removing parameters, etc. The system should support all of these operations.
- **An on-disk format for changes.** This format should only store the information required to reconstruct the change and should support all update types.
- **Restoring the previous state of a checkpoint.** The system should be able to rapidly restore the state of the model at a certain point of history.
- **Communicate/receive history.** There should be a way for me to communicate my local history to a model to someone/someplace else, as well as bring in the changes to the model someone else has proposed.

- **Track changes to a checkpoint.** There are many ways to modify a checkpoint, e.g. updating all parameters, updating a subset of parameters, adding/removing parameters, etc. The system should support all of these operations.
- **An on-disk format for changes.** This format should only store the information required to reconstruct the change and should support all update types.
- **Restoring the previous state of a checkpoint.** The system should be able to rapidly restore the state of the model at a certain point of history.
- **Communicate/receive history.** There should be a way for me to communicate my local history to a model to someone/someplace else, as well as bring in the changes to the model someone else has proposed.
- **Identifying and possibly resolving merge conflicts.** When there are conflicting changes created in parallel, the system should fail or (ideally, eventually) try to resolve them.

- Track changes to a checkpoint. There are many ways to modify a checkpoint, e.g. updating all parameters, updating a subset of parameters, adding/removing parameters, etc. The system should support all of these operations.
- **An on-disk format for changes.** This format should only store the information required to reconstruct the change and should support all update types.
- **Restoring the previous state of a checkpoint.** The system should be able to rapidly restore the state of the model at a certain point of history.
- **Communicate/receive history.** There should be a way for me to communicate my local history to a model to someone/someplace else, as well as bring in the changes to the model someone else has proposed.
- **Identifying and possibly resolving merge conflicts.** When there are conflicting changes created in parallel, the system should fail or (ideally, eventually) try to resolve them.
- Easy to go from training → version control system. We don't expect contributors to directly create patches; we expect them to use version control-aware training code that needs to interact with the system to create patches, etc.

□ r-three / git-cml Public	State Pins -	• Unwatch 5	▼ <sup>Q</sup> Fork 2 ▼ ☆ Star 0 ▼
<> Code     Issues 8	) Pull requests 🖓 Discussions 🕑 Actions 🗄	] Projects 🖾 W	Viki 🕛 Security 🗠 Insights 🛛 …
ີະ° main → ເγ 2 branches	So to file Add file	<> Code •	About 贷
Image: Inclusion of the second state of the second stat			
.github/workflows	Add auto-linting	17 days ago	대 Readme 최 Apache-2.0 license ☆ 0 stars
🖿 bin	Allow users to clone properly from a remote.	18 seconds ago	
examples	Git cml init (#32)	12 days ago	
📄 git_cml	Allow users to clone properly from a remote.	18 seconds ago	S watching
🗋 .gitignore	Allow bin directory	17 days ago	父 2 forks
LICENSE.md	Add license (#40)	4 days ago	Releases
🖺 README.md	Create README.md (#39)	9 days ago	
🗋 setup.py	Add torch to requirements (#41)	10 days ago	No releases published Create a new release

How can we enable collaborative and continual development of machine learning models?

Let's try it by focusing on a specific architecture.

![](_page_44_Figure_0.jpeg)

![](_page_45_Figure_0.jpeg)

![](_page_46_Figure_0.jpeg)

![](_page_47_Figure_0.jpeg)

![](_page_48_Figure_0.jpeg)

![](_page_49_Figure_0.jpeg)

![](_page_50_Figure_0.jpeg)

![](_page_51_Figure_0.jpeg)

- **It would probably work.** Parameter-efficient fine-tuning methods work well. We can probably get the learned routing to work reasonably well.

- **It would probably work.** Parameter-efficient fine-tuning methods work well. We can probably get the learned routing to work reasonably well.
- **It would be useful.** The model would be able to perform tons of tasks out-of-the-box. This would make it convenient and compelling to use.

- **It would probably work.** Parameter-efficient fine-tuning methods work well. We can probably get the learned routing to work reasonably well.
- **It would be useful.** The model would be able to perform tons of tasks out-of-the-box. This would make it convenient and compelling to use.
- Most version control operations are trivial. People can add or update specific adapters, which results in cheap-to-communicate updates that are trivial to merge. Testing would be easy too just evaluate performance on the task for the updated adapter. Updating/adding adapters is probably how most people would contribute.

- **It would probably work.** Parameter-efficient fine-tuning methods work well. We can probably get the learned routing to work reasonably well.
- **It would be useful.** The model would be able to perform tons of tasks out-of-the-box. This would make it convenient and compelling to use.
- Most version control operations are trivial. People can add or update specific adapters, which results in cheap-to-communicate updates that are trivial to merge. Testing would be easy too just evaluate performance on the task for the updated adapter. Updating/adding adapters is probably how most people would contribute.
- **Could occasionally update the router or backbone.** These can be seen as minor/major version updates...

- **It would probably work.** Parameter-efficient fine-tuning methods work well. We can probably get the learned routing to work reasonably well.
- **It would be useful.** The model would be able to perform tons of tasks out-of-the-box. This would make it convenient and compelling to use.
- Most version control operations are trivial. People can add or update specific adapters, which results in cheap-to-communicate updates that are trivial to merge. Testing would be easy too just evaluate performance on the task for the updated adapter. Updating/adding adapters is probably how most people would contribute.
- **Could occasionally update the router or backbone.** These can be seen as minor/major version updates...
- **Could consider adapters beyond task-level.** For example, we could say each task has some particular domain and language, and have separate domain and language adapter sets.

- **It would probably work.** Parameter-efficient fine-tuning methods work well. We can probably get the learned routing to work reasonably well.
- **It would be useful.** The model would be able to perform tons of tasks out-of-the-box. This would make it convenient and compelling to use.
- Most version control operations are trivial. People can add or update specific adapters, which results in cheap-to-communicate updates that are trivial to merge. Testing would be easy too just evaluate performance on the task for the updated adapter. Updating/adding adapters is probably how most people would contribute.
- **Could occasionally update the router or backbone.** These can be seen as minor/major version updates...
- **Could consider adapters beyond task-level.** For example, we could say each task has some particular domain and language, and have separate domain and language adapter sets.
- **Some precedent in AdapterHub.** This model could be seen as a next-generation version; here, adapters ship with the model, routing can be learned, and there is a more principled way of tracking changes.

# Join us!

# https://bit.ly/cccml-community

# Give me feedback:

https://bit.ly/colin-talk-feedback