

A MACHINE LEARNING RESEARCH TOOLKIT

Anderson de Andrade

School of Engineering Science
Simon Fraser University
Burnaby, BC, Canada

August 16, 2025

- Rationale
- Philosophy
- Architecture breakdown
- Component dive-in:
 1. Code
 2. Artifacts
 3. Deployment

1. Before working on theory, might be best to produce empirical results - requires **development speed** and **correctness**
2. In absence of theoretical proof, **correctness** of empirical results takes precedence
3. Dynamic allocation of computing resources essential for **quick turnaround**
4. Recording, labeling and management of experiment artifacts is essential for result **validity** and **reproducibility**
5. Most of your research is **extensible**, so should your code

1. Simple over easy¹

- One role, one task, one concept, one dimension
- But **not**: one instance, one operation, one function, one script
- Often means making more things, not fewer
- Does **not** mean: convenient, available, familiar, at hand, succinctly described
- *"Simplicity is prerequisite for reliability"* - Edsger Dijkstra
- *"Simplicity is the ultimate sophistication"* - Leonardo da Vinci

¹Rich Hickey. Simple Made Easy. StrangeLoop 2011

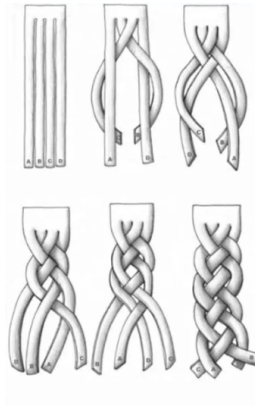


Figure 1: Which one would you rather work with?²

²Rich Hickey. Simple Made Easy. StrangeLoop 2011

2. Decoupling of concerns is imperative

- Ease of understanding
- Ease of change
- Flexibility: pick what you need/like
- Modularity and organization do not imply it but enable it
- Libraries over frameworks: remain in control
- ~~Convention over configuration~~ Explicit is better than implicit



(a) Only need to know what (interface)



(b) Need to also know how^a

^aRich Hickey. Simple Made Easy. StrangeLoop 2011

3. Justified abstractions

- To draw away, not to hide
- About separating the **what** and **how**, also who, when, and where
- ~~Don't repeat yourself~~ Avoid hasty abstractions
- Keep current research your priority
- How much time are we saving?
- Opinion will decrease over time: perspective, not skill

- Extend philosophy to the tools/libraries you use: choose accordingly
- Better to reimplement/port than bring hairball to codebase
- Be able to reason about entire system

ARCHITECTURE BREAKDOWN

Type	Constructs	Tools ³
Code	Versioning	GitHub
	Project template	Cookiecutter
	Dependencies, environment, packaging	Poetry, PyEnv
Code	Model development	PyTorch, Numba
	Training & inference boilerplate	PyTorch Lightning
	Dataset development	PyTorch Data Utilities
	Data augmentations	TorchVision, SFU Torch Library
	Argument parsing	Defopt

³All tools instrumented in project template

ARCHITECTURE BREAKDOWN

Type	Constructs	Tools ⁴
Code	Formatting Type checking Tests	ruff basedpyright PyTest
Artifacts	Dataset storage Experiment artifacts Notifications	MinIO, SFU Torch Library MLFlow Tracking, SFU Torch Library Slack, SFU Torch Library
Deployment	Containers Job scheduling Credentials, configuration	Docker Images, MLFlow Projects Slurm, Apptainer & Kubernetes Ansible, Direnv

⁴All tools instrumented in project template

DIVE-IN: VERSIONING (CODE)

- Uncluttered: **avoid** commented blocks, extra files, unused implementations
- Rely on **meaningful** code versioning
- Experiments run on specific version
- OK to rewrite history outside main branch: **git-aap**
- Codebase is not your agenda: LogSeq, Obsidian
- Cannot memorize the CLI? Lazygit, VSCode, GitHub Desktop

- Provides:
 - Minimal project configuration using collected information
 - *Make* tasks for deployment, linting, and testing
 - Job scheduling scripts
 - Small project scaffold
- Duplication allows project independence
- Run `make project`

DIVE-IN: DEPENDENCIES, ENVIRONMENT, PACKAGING (CODE)

- *Poetry lock* enables reproducibility
- *PyEnv* manages different *Python* versions, installed under user
- Include all dependencies that you explicitly import
- Be specific about version unless used mainly by dependencies
- Treat *Python* version as a dependency
- Initialize environment `make init`
- Activate environment `$(poetry env activate)`
- Run command within environment `poetry run`
- Build package `make build`

- *PyTorch's* **Module** interface enables decoupling and abstractions
- Express core logic in **simple** functions
- Group functions with interface implementations
- Keep interfaces pure
- Enable *Accelerated Linear Algebra* (XLA) with `torch.compile()`
- Custom operators with **Function**⁵ and `torch.library`
- Low-level implementations with CUDA support using *Numba*⁶

⁵`torch.autograd.Function`

⁶See numba.pydata.org

DIVE-IN: TRAINING & INFERENCE BOILERPLATE (CODE)

- *PyTorch Lightning*⁷ is a minimalist **framework**
- Provides a high-level interface for *PyTorch* models: **LightningModule**
- **Trainer** provides:
 - Gradient management
 - Dataloader management: device placement
 - Metric accumulation: supports torchmetrics⁸
 - **Callbacks**
- Callbacks are main way to manage functionality

⁷See lightning.ai/docs/pytorch/stable

⁸See lightning.ai/docs/torchmetrics/stable

- *PyTorch's Dataset*⁹ and **IterableDataset** interfaces are simple
- Manipulations: stacking, concatenation, chaining, subsets
- Different sampling methods, customizable via **Sampler**
- **FileFetcher**¹⁰ interface implemented for ZIP, TAR and file systems
- **DataLoader**: multiprocessing data loading, prefetching and memory pinning
- Control over batch collation with overloaded defaults

⁹`torch.utils.data`

¹⁰`sfu_torch_lib.file_fetcher`

- *TorchVision's* **Compose**¹¹ for transformation pipelines
- Use tensors instead of *PIL* images
- Use `torch.uint8` dtype
- Use *SFU Torch Lib's* **ComposeTree**¹² for fine-grained group transformations

¹¹`torchvision.transforms.v2`

¹²`sfu_torch_lib.processing`

- **defopt**:¹³ lightweight, minimal effort
- Call functions without having to recreate arguments object
- Allows configuration to be explicit
- Only documented functions can become abstractions

¹³See defopt.readthedocs.io

DIVE-IN: FORMATTING (CODE)

- Readability and consistency across projects
- `ruff`¹⁴ written in Rust (fast)
- 800+ **configurable** rules
- `pyproject.toml` support
- Automatic error correction
- Run `make format`
- Keep consistency between editor and CLI

¹⁴See github.com/astral-sh/ruff

DIVE-IN: TYPE CHECKING (CODE)

- HPC is a strong case for static typing
- `basedpyright`¹⁵ forks `pyright`
- Adds `pylance`¹⁶ for completion, quick info, members
- `pyproject.toml` support
- Assertions enforce type
- Avoid reusing variables, simple
- `# type: ignore`
- Run `make type-check`
- Keep consistency between editor and CLI

¹⁵See docs.basedpyright.com

¹⁶See marketplace.visualstudio.com/items?itemName=ms-python.vscode-pylance

- No more than guardrails
- HPC has a strong case for it
- Unit test vs. functional (integration) test
- Useful assertions in `torch.testing`¹⁷ and `numpy.testing`¹⁸
- Run `make test`

¹⁷pytorch.org/docs/stable/testing.html

¹⁸numpy.org/doc/stable/reference/routines.testing.html

DIVE-IN: DATASET STORAGE (ARTIFACTS)

- Centralized, decoupled and ubiquitous
- Available via Internet
- S3-compatible for extensive API support¹⁹
- Two-level caching via `sfu_torch_lib.io`²⁰
- `aws s3` for management
- GUI `minio.multimedialabsfu.ca`

¹⁹On Python, use boto3

²⁰See `sfu_torch_lib.io`

DIVE-IN: EXPERIMENT ARTIFACTS (ARTIFACTS)

- MLFlow Tracking focuses on experimentation
- Parameters are considered immutable
- Datasets are often immutable, only URL as parameter
- Runs linked to Git version
- Instrument via logger, checkpoint callback and entry function annotator
- GUI `mlflow.multimedialabsfu.ca`
- `mlflow-skinny`²¹ API allows to retrieve data programmatically
- Uses our Minio cluster for model storage

²¹See pypi.org/project/mlflow-skinny

DIVE-IN: NOTIFICATIONS (ARTIFACTS)

- Slack notification on script: start, end and error
- Collects arguments from function call
- Collects and reports exception messages
- Implemented via entry function annotator²²

²²See `sfu_torch_lib.slack`

DIVE-IN: CONTAINERS (DEPLOYMENT)

- Docker images are OCI-compliant
- Compatible with: Apptainer, Kubernetes
- Embeds project's Python and dependencies
- MLFlow Project²³ for entry points
- **MLProject** included in scaffold
- Build `make training-image`
- Shared repository on DockerHub²⁴ with project name as label

²³See mlflow.org/docs/latest/projects

²⁴See hub.docker.com/repositories/multimedialabsfu

DIVE-IN: JOB SCHEDULING (DEPLOYMENT)

- Project commands²⁵ support: Slurm, Kubernetes and Apptainer
- `slurm-{launch,status,log,stop,cancel}`
- Slurm scripts support periodic restarts
- Support for argument sets
- Relaunch with `slurm-launch --run-id`

²⁵See repository

DIVE-IN: CREDENTIALS & CONFIGURATION (DEPLOYMENT)

- Credentials stored in encrypted file
- Decryption key stored in OS keychain
- All credentials and configuration as environment variables
- Loaded in shell when changing to project directory: `direnv allow`
- Available in job's environment

- Fork github.com/multimedialabsfu/research
- Follow instructions on **README.md** to install required software
- Update cookiecutter before creating a new project