

Attempt on Cerebras:

```
File "<frozen importlib._bootstrap>", line 219, in _call_with_frames_removed
File "<frozen importlib._bootstrap>", line 1014, in _gcd_import
File "<frozen importlib._bootstrap>", line 991, in _find_and_load
File "<frozen importlib._bootstrap>", line 961, in _find_and_load_unlocked
File "<frozen importlib._bootstrap>", line 219, in _call_with_frames_removed
File "<frozen importlib._bootstrap>", line 1014, in _gcd_import
File "<frozen importlib._bootstrap>", line 991, in _find_and_load
File "<frozen importlib._bootstrap>", line 961, in _find_and_load_unlocked
File "<frozen importlib._bootstrap>", line 219, in _call_with_frames_removed
File "<frozen importlib._bootstrap>", line 1014, in _gcd_import
File "<frozen importlib._bootstrap>", line 991, in _find_and_load
File "<frozen importlib._bootstrap>", line 973, in _find_and_load_unlocked
ModuleNotFoundError: No module named 'cerebras.modelzoo.models.nlp.bert.venv_cerebras_pttorch==2'

During handling of the above exception, another exception occurred:

Traceback (most recent call last):
  File "run.py", line 41, in <module>
    main()
  File "run.py", line 32, in main
    main()
  File ".../cerebras/modelzoo/common/run_utils.py", line 115, in main
    return run_with_params(
  File ".../cerebras/modelzoo/common/run_utils.py", line 182, in run_with_params
    return run_trainer(
  File ".../cerebras/modelzoo/trainer/utils.py", line 85, in run_trainer
    trainer = configure_trainer_from_params(params, model_fn)
  File ".../cerebras/modelzoo/trainer/utils.py", line 281, in configure_trainer_from_params
    _validate_trainer_params(trainer_params, model_fn_or_name)
  File ".../cerebras/modelzoo/trainer/utils.py", line 1128, in _validate_trainer_params
    if (config_class := registry.get_config_class(model_name)) is None:
  File ".../cerebras/modelzoo/common/registry.py", line 360, in get_config_class
    cls._import_modules()
  File ".../cerebras/modelzoo/common/registry.py", line 96, in _import_modules
    cls._import_modules_for_registry(
  File ".../cerebras/modelzoo/common/registry.py", line 69, in _import_modules_for_registry
    raise Exception("Registry Import Failure: {}".format(ex))
Exception: Registry Import Failure: No module named 'cerebras.modelzoo.models.nlp.bert.venv_cerebras_pttorch==2'
(venv_cerebras_pt) [ed Doyle@cer-login-03 bert]$
```

Attempt on Groq:

Gave “conda command not found” after installing miniconda, and couldn’t get around this.

Attempt with my friend on her computer with Cerebras: (we did the same thing and hers worked while mine did not...)

```
andreadiaz@cer-login-03:~/F x + v
2024-11-21 01:17:16,757 INFO: Exploring data layouts
2024-11-21 01:18:51,601 INFO: Optimizing memory usage
2024-11-21 01:20:30,064 INFO: Gradient accumulation trying micro batch size 64...
2024-11-21 01:20:44,030 INFO: Exploring floorplans
2024-11-21 01:21:00,890 INFO: Exploring data layouts
2024-11-21 01:21:53,115 INFO: Optimizing memory usage
2024-11-21 01:23:01,431 INFO: Gradient accumulation trying micro batch size 512...
2024-11-21 01:23:16,430 INFO: Exploring floorplans
2024-11-21 01:23:34,281 INFO: Exploring data layouts
2024-11-21 01:25:33,879 INFO: Optimizing memory usage
2024-11-21 01:26:46,090 INFO: Gradient accumulation trying full batch size 1024...
2024-11-21 01:26:59,925 INFO: Exploring floorplans
2024-11-21 01:27:09,069 INFO: Exploring data layouts
2024-11-21 01:29:19,519 INFO: Optimizing memory usage
2024-11-21 01:31:00,235 INFO: Gradient accumulation showed a benefit
2024-11-21 01:31:00,839 INFO: Post-layout optimizations for <batch=1024, lanes=11>...
2024-11-21 01:31:00,841 INFO: Post-layout optimizations for <batch=1024, lanes=9>...
2024-11-21 01:31:00,841 INFO: Post-layout optimizations for <microbatch=256, lanes=3>...
2024-11-21 01:31:00,843 INFO: Post-layout optimizations for <batch=1024, lanes=10>...
2024-11-21 01:31:00,848 INFO: Post-layout optimizations for <microbatch=512, lanes=6>...
2024-11-21 01:31:21,167 INFO: Allocating buffers for <microbatch=256, lanes=3>...
2024-11-21 01:31:22,666 INFO: Allocating buffers for <microbatch=512, lanes=6>...
2024-11-21 01:31:23,856 INFO: Allocating buffers for <batch=1024, lanes=10>...
2024-11-21 01:31:24,192 INFO: Allocating buffers for <batch=1024, lanes=9>...
2024-11-21 01:31:25,101 INFO: Code generation for <microbatch=256, lanes=3>...
2024-11-21 01:31:25,455 INFO: Allocating buffers for <batch=1024, lanes=11>...
2024-11-21 01:31:26,809 INFO: Code generation for <microbatch=512, lanes=6>...
2024-11-21 01:31:28,121 INFO: Code generation for <batch=1024, lanes=10>...
2024-11-21 01:31:28,436 INFO: Code generation for <batch=1024, lanes=9>...
2024-11-21 01:31:30,004 INFO: Code generation for <batch=1024, lanes=11>...
2024-11-21 01:31:57,927 INFO: Compiling image...
2024-11-21 01:31:58,043 INFO: Compiling kernels
2024-11-21 01:34:17,011 INFO: Compiling final image
2024-11-21 01:37:22,152 INFO: Compile artifacts successfully written to remote compile directory. Compile hash is: cs_1844636080405
7867627
2024-11-21 01:37:22,337 INFO: Compile was successful!
```

Theory Homework: (the solutions were given, so I paraphrased the notes to aid my learning)

1. What are the key architectural features that make these systems suitable for AI workloads?

AI accelerators enable efficient processing of intricate, memory intensive operations. These are specially designed pieces of hardware (highly engineered layout and connections of transistors on silicon chips) that enable fast tensor and matrix operations. In high dimensional problems memory is a key issue, and the architectures are designed with memory in mind, where there is intricate design of on-chip memory and off-chip memory storage logistics. The systems have various different centers of work called cores, processing units, tiles, etc depending on the system. These interconnected units enable the parallel processing of data by splitting up tasks across the various units so that you can accelerate your training and calculations.

2. Identify the primary differences between these AI accelerator systems in terms of their architecture and programming models.

Sambanovas Reconfigurable Dataflow Unit (RDU) is best for large amounts of data and has multiple tiers of memory storage allowing for creative manipulation of the pathways for processing. Cerebras is based on processing elements that individually handle and store their data allowing for fast parallelization with high scalability. Graphcore has tiles that similarly store and handle their data. The processing breaks tasks into computation and communication steps, instead of continual memory saving. Groq's architecture helps to streamline procedures and minimize over use of memory features to enable consistent and efficient deterministic tasks.

3. Based on hands-on sessions, describe a typical workflow for refactoring an AI model to run on one of ALCF's AI testbeds. What tools or software stacks are typically used in this process?

I did most of my work on Cerebras, so I will discuss the process for that system. All details are found here: <https://docs.alcf.anl.gov/ai-testbed/getting-started/> with details in the side menu for each system. Cerebras operates through communication with the Python package `modelzoo.common.pytorch.run_utils`. Models from PyTorch are used and can be trained and optimized in the system. To do work, you must create a PyTorch virtual environment for Cerebras and clone the modelzoo framework into the environment for use. You then must ensure all directories are edited to work with your personal account. Once all of the details are set you can run the job.

4. Give an example of a project that would benefit from AI accelerators and why?

In my lab we use AI for molecule discovery. We built a massive dataset of battery electrolyte compositions and their performance data to use with a neural network model to predict what new molecules would be optimal as battery electrolytes. In this process we talked a good deal about using DFT supplemented data. I think the AI accelerators could be helpful for parallelizing the DFT calculations to supplement the large data set in the data curating step. I'm sure it would also enable advanced ML techniques beyond basic NNs.