

## General:

- Spike data hierarchy: multiple physical *probes* -> multiple physical *channels* (electrodes)  
-> multiple informational/representational/calculated *units* (neurons)
- Longleaf settings:
  - GPU regular instead of FULL
  - A100
- Additional job submission arguments:
  - --mem=64G (up to 118)
- Jupyter
- \*/proj/STOR/pipiras/Neuropixel\*
- Have to train per subject (probably cross session)
- Zinodo macak cursor data (neurolink) for more data
- Environment:
  - Metadata shows pipeline using Python 3.9.7
  - Python 3 (ipykernel): Using python 3.11, torch 2.4.1+cu121 cuda: 12.1,
    - If using PyG need a PyG build that matches Torch 2.4.1/cu121
  - From “module avail”, “python/3.9.6” is an option I can switch to, but should probably stick to newer kernel env bc backwards compatibility with old torch versions will probably cause bigger issues with newer CUDA driver.

## Todo:

- Ensure data is not corrupted, remove simplistic filters like in data\_splitter->DataSplitter->\_split\_and\_reshape to test for corruption, and fix if needed
- Try to run model predictions .ipynb, look at what is happening from the highest level pipeline perspective and play around with it try to train and test models and see output, reorganize it in my own directory if needed, then deep dive into models.

## Files Types:

- **.pkl** — Python pickles holding model-ready arrays/DataFrames.
- **.pt** — PyTorch tensors (not models) such as the VAE latent encodings of each stimulus image used for feature-space comparisons.
- **.pth** — PyTorch model checkpoints (parameters) for trained networks (e.g., Transformer→LSTM heads that predict spikes/firing).

- **.nwb** — Allen/Neurodata Without Borders session files with raw data and metadata. Contains full unprocessed session.
- **.csv** — Flat tables for metadata and small numeric summaries (e.g., units.csv, channels.csv, sessions.csv) used for joins/filtering.
- **.json** — Manifests/config/settings and SDK cache info (e.g., ecephys\_project\_manifest.json, session maps, preprocessing configs). Parameters for data fetching from SDK.

**Boc:** folder with unused .nwb file, calcium data (not electrical probes)

## Data\_processors: <-

- **Pull\_and\_process\_data.py:** Extract target session info from AllenSDK API, process the spike information into a 2D Pandas dataframe that counts every unit's spikes for every time bin within a frame exposure window, for every frame exposure window. Exports the raw spike counts as a .pkl file to "output" folder.
  - **print\_memory\_usage:** diagnostics helper that prints the current Python process's RAM use.
  - **print\_gpu\_memory:** diagnostics helper that prints the current Python process's GPU RAM use and cached space in reserve.
  - **NeuronDataset:** turns "per-neuron spike time series + shared timing info" into on-demand binned spike vectors, so you can iterate or use a DataLoader(num\_workers=..., batch\_size=...) to parallelize/stream the per-neuron binning step without holding all processed neurons in memory at once.
    - **Process\_neuron\_wrapper:** Wrapper function to unpack arguments and call process\_neuron.
    - **Process\_neuron (modified):** Outputs 1-D tensor vector that corresponds to a unit's spike count for each time ordered bin in which an image was shown.
  - **Create\_directory\_and\_manifest:** Create the output directory and check or create the manifest.json file. Parameters: {directory\_name (str): Name of the directory to be created. Default is 'output'.} Returns: {tuple: A tuple containing the path to the output directory and the path to the manifest.json file.}
  - **Create\_cache\_get\_session\_table:** Create an instance of the EcephysProjectCache class and get the session table. Parameters: {manifest\_path (str): Path to the manifest.json file.} Returns: {tuple: A tuple containing the cache object and the session table.}

- **Pick\_session\_and\_pull\_data (Modified):** Returns “spike\_times” dictionary with unit label as key and the times it spikes as the value (list), and “stimulus\_table”, where each image presentation is a row and 2 columns are “start\_time” and “stop\_time” for when that image is shown. Prints additional attributes/methods of the session object.
  - Also returns: session filtered by the following parameters: (isi\_violations\_maximum=np.inf, amplitude\_cutoff\_maximum=np.inf, presence\_ratio\_minimum=-np.inf,) (the same session used to get spike\_times and stimulus\_table)
- **Filter\_valid\_spike\_times:** Filter the valid spike times using invalid times from the session object, converting the “spike\_times” dic to “valid\_spike\_times”
- **Get\_stimulus\_table:** Retrieve the stimulus table for a given stimulus type.
- **Calculate\_bins:** thin wrapper returning 4 items
- **Process\_all\_neurons (modified):** Process unit spike times into model ready binned spike count vectors for all neurons, returns CPU-stored list of 1-D tensor for each unit. Processes in parallel using batches of 1000.
  - If you intend to keep results on GPU for downstream ops, don’t wrap with torch.tensor(...) at return; just return spike\_matrix. If you want an actual 2-D tensor later, call torch.stack(spike\_matrix) when needed (mind memory).
- **Create\_and\_prepare\_spike\_dataframe (modified):** Create a spike DataFrame and prepare it by adding a frame columnn.
 

Parameters:

  - spike\_matrix (torch.Tensor): The spike matrix containing processed spike times for all neurons. Is not directly the output of Process\_all\_neuronsbc that is a list.
    - Modified: spike\_matrix (list): A list of 1000-neuron spike matrix batches containing processed spike times for all neurons.
  - spike\_times (dict): A dictionary containing the original spike times.
  - stimulus\_table (pd.DataFrame): The stimulus table containing stimulus-related data.
  - timesteps\_per\_frame (int): Number of timesteps per frame for binning.

Returns:

  - pd.DataFrame: A prepared DataFrame containing spike data and a frame column.
- **Save\_and\_count\_spike\_dataframe:** Save the spike DataFrame to a pickle file in the specified directory and count the number of NaN values.
 

Parameters:

  - spike\_df (pd.DataFrame): The spike DataFrame to be saved.

- session\_number (int): The session number.
- output\_dir (str): The directory where the DataFrame should be saved.
- timesteps\_per\_frame (int): Number of timesteps per frame for binning.

Returns:

- str: The name of the saved file.
- **Normalize\_firing\_rates:** Normalize the firing rates by calculating z-scores.
- **Filter\_and\_save\_neurons:** filter neurons based on firing rate thresholds, check for NaN values, and save the filtered data.

Parameters:

- normalized\_firing\_rates (pd.DataFrame): DataFrame containing normalized firing rates.
- highest\_value (float): Upper z-score threshold for filtering neurons.
- lowest\_value (float): Lower z-score threshold for filtering neurons.
- session\_number (int, optional): The session number, required if saving the data.
- output\_dir (str, optional): The directory where the DataFrame should be saved, required if saving the data.

Returns:

- pd.DataFrame: DataFrame containing filtered neurons.
- **Process\_and\_save\_data: Outdated**
- **get\_session\_ids:** Get all available session IDs from the cache.
- **Process\_session\_batch: Outdated**
- **Process\_all\_sessions\_in\_batches: Outdated**
- **Master\_function (modified):** Master function to execute the entire workflow from data extraction to saving. **Normalized firing rates functionality commented out.**

Parameters:

- session\_number (int): The session number.
- output\_dir (str): The directory where data should be saved. Default is 'output'.
- timesteps\_per\_frame (int, optional): Number of timesteps per frame for binning. Default is 10.

Returns:

- str: The name of the saved filtered data file.
- **Load\_processed\_data.py:** takes raw spike time .pkl file and normalizes and filters out units with abnormal spike counts (above 100 or below 0), as well as frames with values -1 (probably unnecessary). Exports cleaned .pkl file to “output” folder.

- **check\_file\_exists (Modified):** Check if the pickle file exists in the specified directory.

Parameters:

- session\_id (int): The ID of the session.
- output\_dir (str): The directory where the pickle file is expected to be found.
- file\_prefix (str): The prefix of the pickle file.

Returns:

- bool: True if the file exists, False otherwise.

- **load\_pickle\_file:** Load a pickle file from the given path.

Parameters:

- pickle\_path (str): Path to the pickle file.

Returns:

- object: The object loaded from the pickle file, or None if loading fails.

- **save\_pickle\_file:** Save data to a pickle file at the given path.

Parameters:

- data (object): The data to save.
- pickle\_path (str): Path to save the pickle file.

- **clean\_avg\_firing\_rates:** Clean the average firing rates in the given DataFrame.

Parameters:

- df (DataFrame): The DataFrame with the firing rates data.
- highest\_value (float): Upper z-score threshold for filtering neurons.
- lowest\_value (float): Lower z-score threshold for filtering neurons.

Returns:

- DataFrame: The cleaned DataFrame.

- **master\_cleaning\_and\_saving(Modified):** Master function to check, clean, and save the firing rate data within an 'output' directory. **Remove line removing rows with -1 frame values? (no reason why this would happen)?**

Parameters:

- session\_id (int): The ID of the session to be processed.
- output\_dir (str): The directory where the pickle file is located.
- original\_pickle\_prefix (str): The prefix of the original pickle file.
- new\_pickle\_prefix (str): The prefix for the new cleaned pickle file.
- highest\_value (float): The threshold to remove neurons with a high z-score.
- lowest\_value (float): The threshold to remove neurons with a z-score never exceeding this value.

- **data\_splitter.py:** It prepares a time-ordered dataframe (last column = class) for sequence models by label-encoding the targets, dropping zero-variance features, splitting (no shuffle) into train/val/test, segmenting into fixed-length sequences, and reshaping to tensors of shape (B, T, N, 1) while using the label from the last timestep of each sequence. It then builds PyTorch DataLoaders for each split and includes a helper to compute a feature correlation matrix.
  - **DataSplitter (class):**
    - **\_\_init\_\_:** initialize X and y attributes using a further processed session pandas dataframe (remove\_0\_columns, \_prepare\_data)
    - **\_prepare\_data:** split session data pandas dataframe into X and y, convert semantic y labels to integers.
    - **Remove\_0\_columns:** remove units that contain zero information, i.e don't vary in their spikes at all (inactive)
    - **\_split\_and\_reshape:** reshape feature pandas dataframe into 3-D (with an added 4th dimension for CNN compatibility). First dimension is each individual sample, i.e each time an image was shown. Each of these is associated with a 2-D matrix of the 2nd and 3rd dimension, which is each time ordered bin and each unit. Label pandas dataframe is reshaped from one label per time bin to one label per sample (with an added 2nd dimension for CNN compatibility)
    - **Potential issue:** "[ :num\_samples \* self.seq\_len]" safety check is flawed because it only solves issues at the tail. If corruption happened in the beginning or middle, this filter will corrupt data.
    - **\_create\_loaders (Modified):** "Master function" of class, all previous helper functions utilized to create cleaned, batched, data loader class attributes for train, test, split (as well as the raw X/y tensor splits. X and y are still pandas dataframes but X\_train etc. are pytorch tensors).
      - **Shuffle batches set to True for train loader to reduce potential overfitting.**
    - **Compute\_correlation\_matrix:** computes the correlation matrix of the units with each other.

## Models:

## Old Files: irrelevant

**Output:** Extracted and Processed data/artifacts used by models

- **Session\_#:** folder containing that session's .nwb file

- **Channels.csv:** contains individual channel ids from the probes and their spatial location
- **Manifest.json:** configurations for fetching data from SDK
- **Probes.csv:** probe id's and technical specifications
- **Session\_details.csv:** LSTM-specific training run hyperparameters and results
  - Session id's do not line up with .nwb or .pkl files in output parent folder
  - Some train\_acc values are > 1 (e.g., 1.12, 1.37) while others look like percentages (92.69) — this suggests a units/logging inconsistency (fraction vs percent). You may want to standardize the logging.
- **Sessions.csv:** largely irrelevant metadata about the session and it's mouse
- **Spike\_trains\_plot.png:** visualization of spike data
- **Spike\_trains\_with\_stimulus\_session\_(session\_id)\_(bins).pkl:**
  - Session pkl's do not line up 1-1 with session .nwb files.
  - your code has hard-coded/changed output\_dir paths at times, so artifacts may have landed in different folders on different runs.
  - Only the 250 ms stimulus frames are included, binned into a set number of intervals, with one vector per interval.
  - Number after session id corresponds to the timesteps\_per\_frame parameter (number of intervals) in the processing code (pull\_and\_process\_data.py / process\_and\_save\_data).
  - Each vector has a “frame” feature (last), and many unit features (activations for each neuron, as ints).
  - Each unit feature's data is a int count of many times the neuron spiked within the time frame of the current bin (row)
  - Rows (bins) time ordered. One training example corresponds to a specific frame's block of time bins during which that frame was shown. A single frame will have multiple training example blocks because they were shown multiple times (roughly 50)
- **Units.csv:** metadata about unit quality and physical characteristics

**Pictures:** high level visualizations of a couple model information flows

## Results:

- **3\_shot\_(model type)\_per\_mouse.csv:** For each image (frame .pkl feature), 3 of its time blocks are randomly selected and used to train the model, instead of the full 50. Model hyperparameter and performance reported.
  - Should be per session, not mouse, one mouse can have multiple sessions
- **highest\_test\_accuracy.png:** Best LSTM accuracy (best model) per mouse session
  - Title slightly misleading one mouse can have multiple sessions

- **Highest\_test\_accuracy\_no\_subs.png:** ^

## Saved\_models:

**Wandb:** “weights and biases” folder to track performance, hyperparameters, env of history of training run. Hyperparameters and weights and biases saved elsewhere so safe to delete, folder is just for training run evaluation

- **Run2024...7j:** Static ST-GAT training run
  - **Files:**
    - **Config.yaml:** hyperparamters and env info of training run
    - **Output.log:** per epoch data
    - **Requirements.txt:** Python environment snapshot: a pinned list of packages (with exact versions) needed to reproduce run
    - **Wandb-summary.json:** Most recent snapshot of performance called at a specific epoch.
  - **Logs:** Debugging information