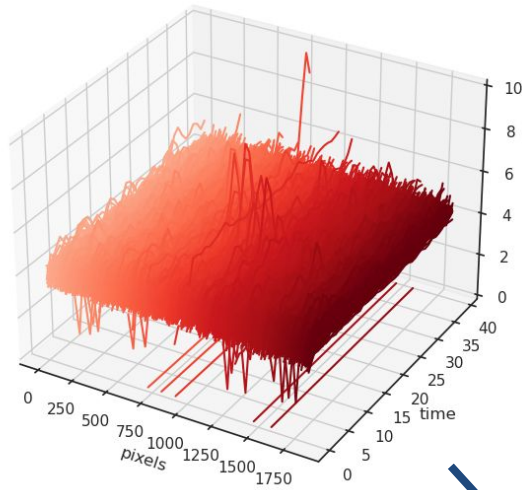


# CTAO Real Time Analysis : investigating porting to GPU

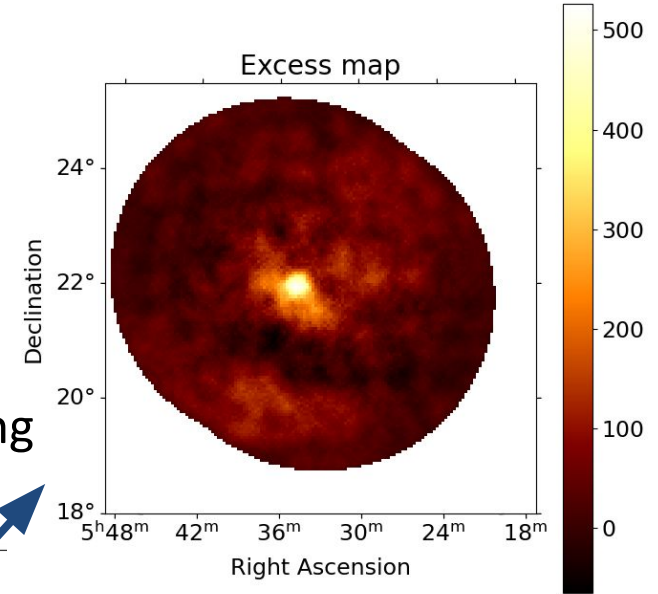
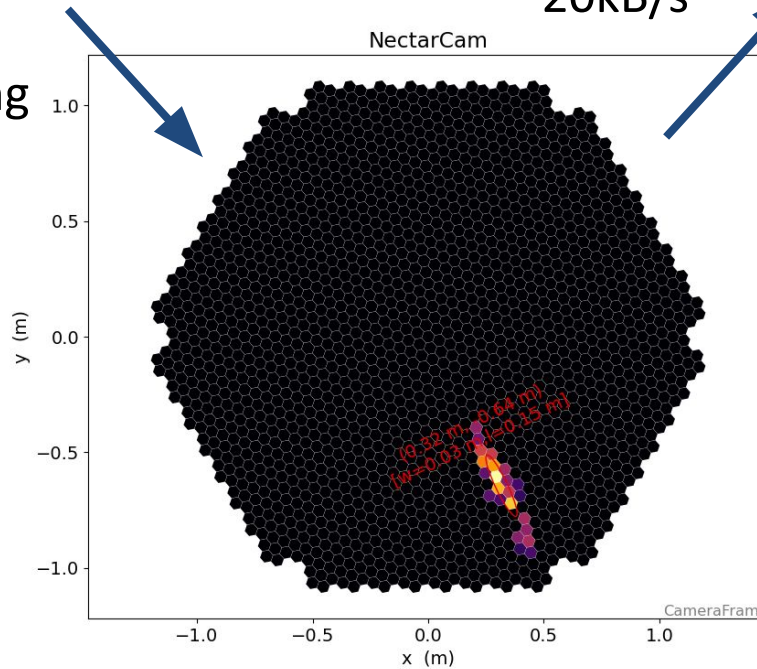
LAPP: Vincent Pollet, Pierre Aubert, Sami Caroff, Thomas Vuillaume  
11/07/2024





Signal processing  
~3GB/s

Machine learning  
~20kB/s



The logo for hiperRTA, featuring the word "hiper" in a red, pixelated font and "RTA" in a bold, black, sans-serif font.The logo for ctapipeline, featuring a circular graphic with a colorful, abstract pattern of lines and dots, followed by the text "ctapipeline" in a bold, black, sans-serif font.

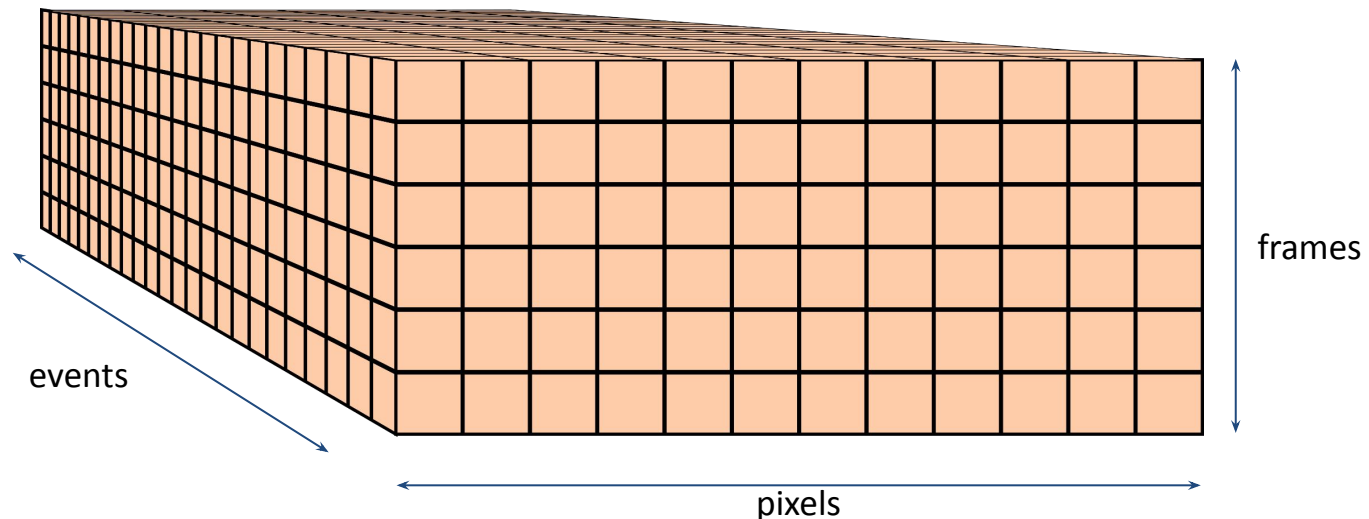
**R&D project: numpy implementation that translates to gpu**

cuNumeric

The logo for PyTorch, featuring a red flame-like shape followed by the text "PyTorch" in a black, sans-serif font.

<https://gitlab.in2p3.fr/CTA-LAPP/rta/pyhiperta/>

- Data:
  - Events:  $\sim 10^7$
  - samples: 40
  - pixels: 1855
- **batch** of few events on CPU
- **batch** of many events on GPU
- Spatial & temporal memory locality  $\rightarrow$  data tensor
- Parallelization  $\rightarrow$  independent events processed in parallel
- vectorization  $\rightarrow$  independent pixel processing per time slice



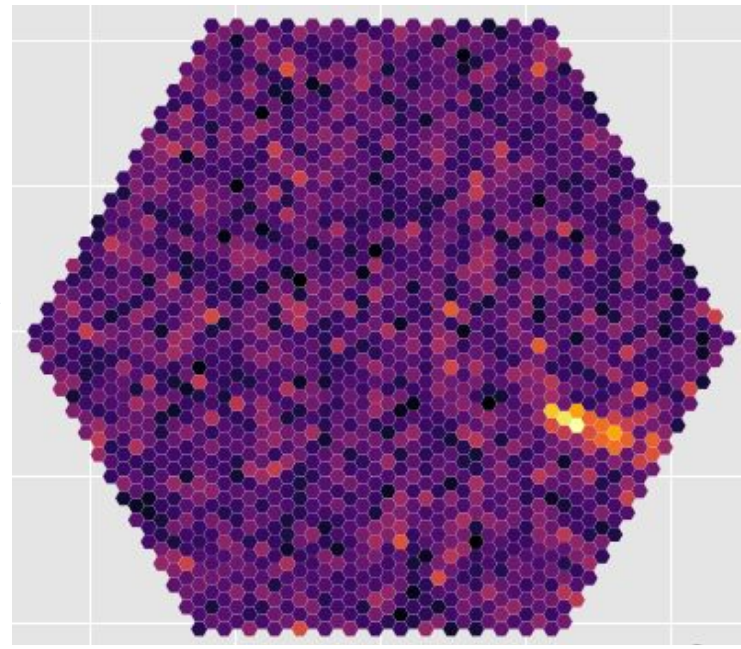
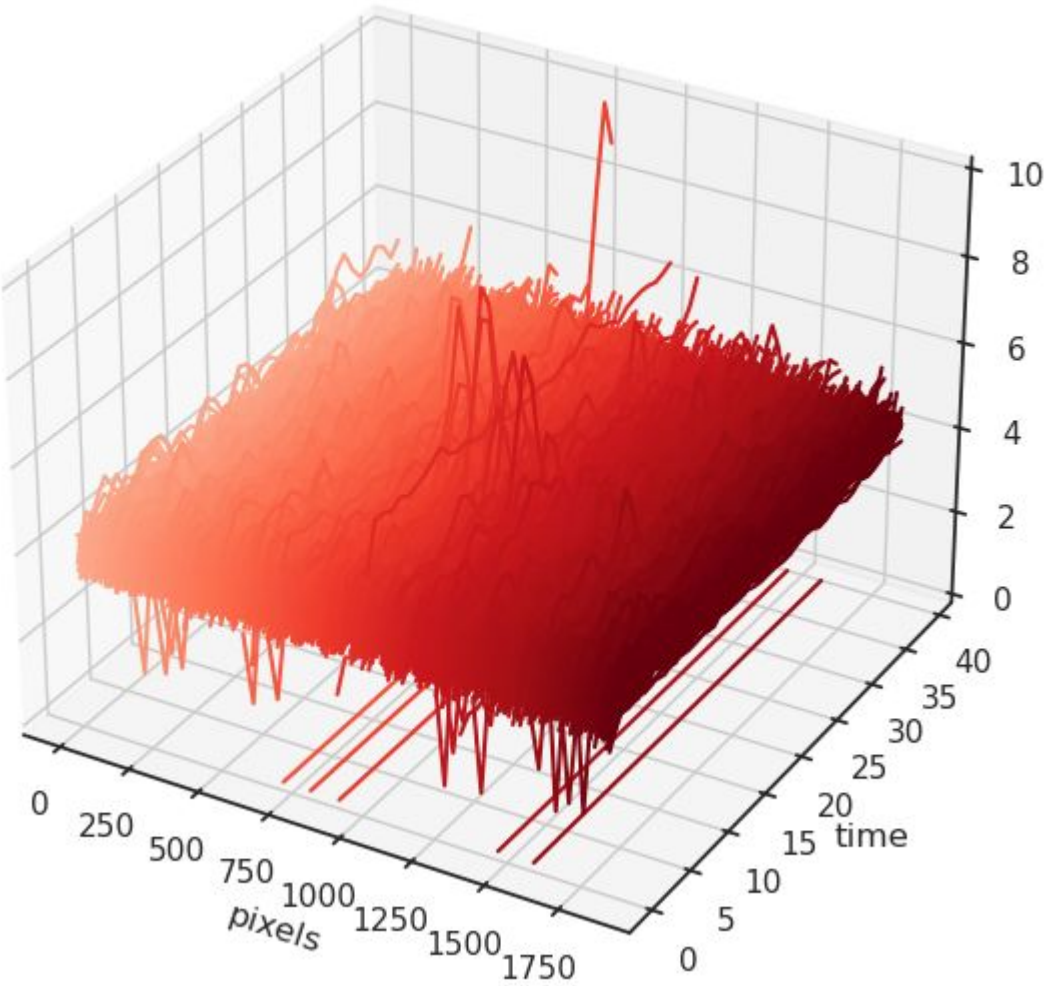
Processing steps:

- calibration
- integration
- cleaning
- hillas parameters

$$w_{e,t,p} = (w_{e,t,p} - ped_{e,p}) \cdot gain_{e,p}$$

```
def calibrate(waveform, gains, pedestals):  
    """  
    waveform  
        Shape: (N_batch, N_frames, N_pixels)  
    gains  
        Shape: (N_batch, N_pixels).  
    pedestals  
        Shape (N_batch, N_pixels)  
    """  
  
    return (waveform - pedestals[..., np.newaxis, :]) * gains[..., np.newaxis, :]
```





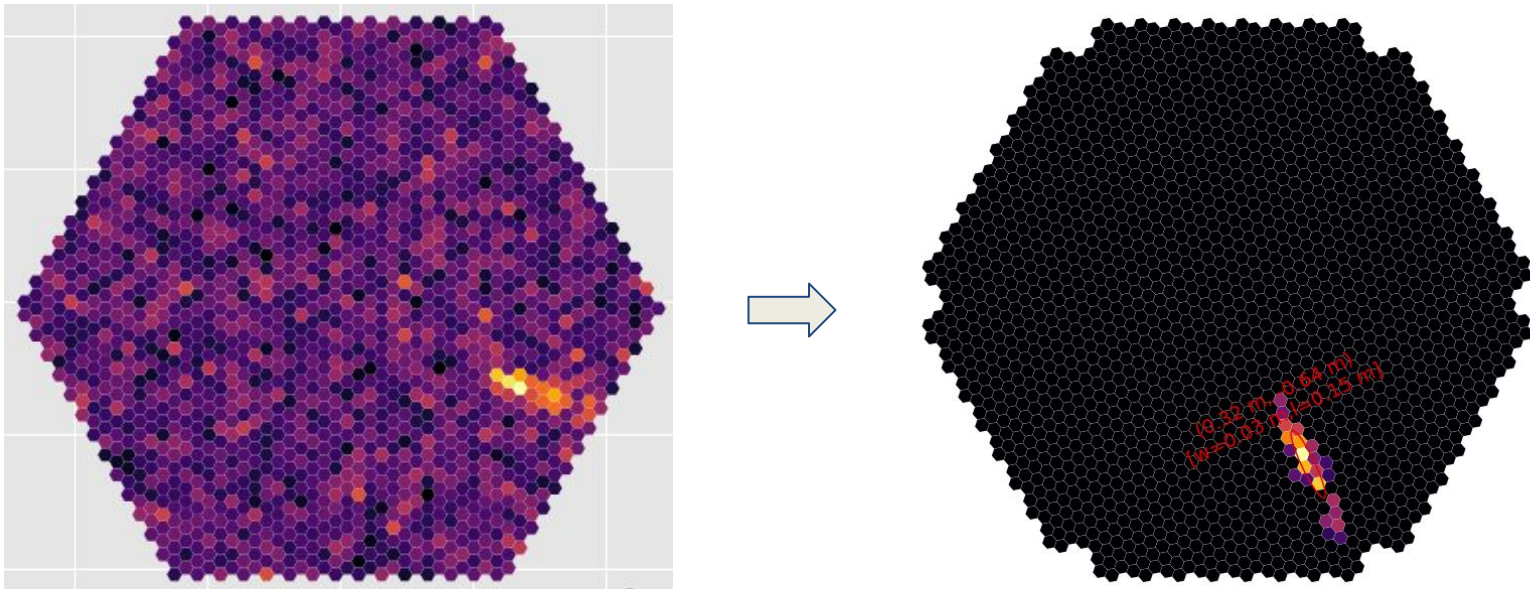


```
def integrate_local_peak(waveforms, nb_frames_before_max, nb_frames_after_max):  
    charge = np.empty((nb_channels, nb_frames))  
    for w in waveforms:  
        for p in waveforms[w]:  
            max_idx = np.argmax(waveforms[w, p])  
            charge[w, p] = waveform(w, p, max_idx-nb_frames_before_max, max_idx+nb_frames_after_max+1,  
                                   p)  
            charge[w, p].sum()  
    return charge
```

Sequential operations

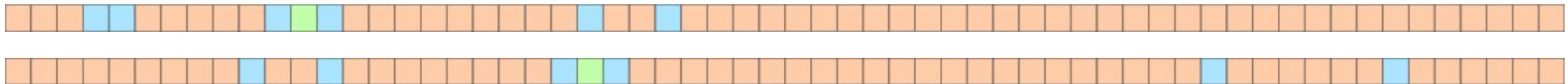
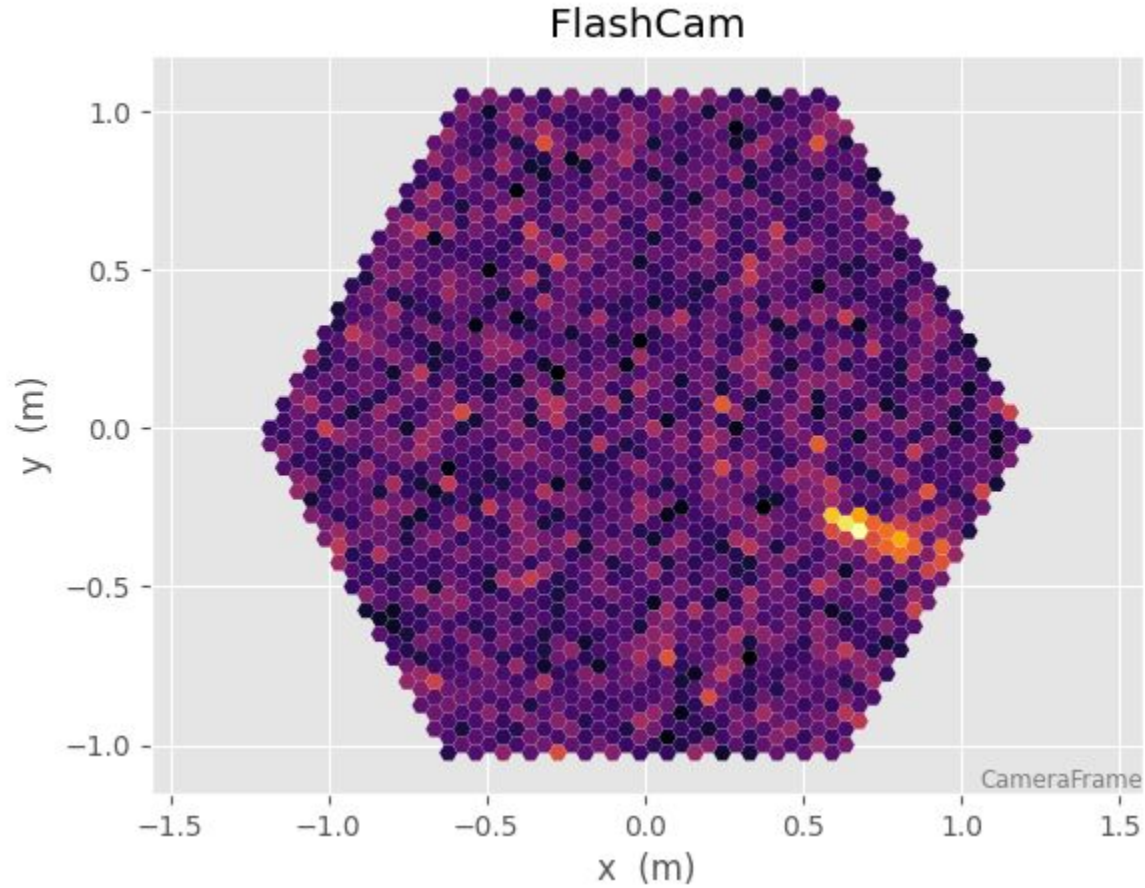
on vectors

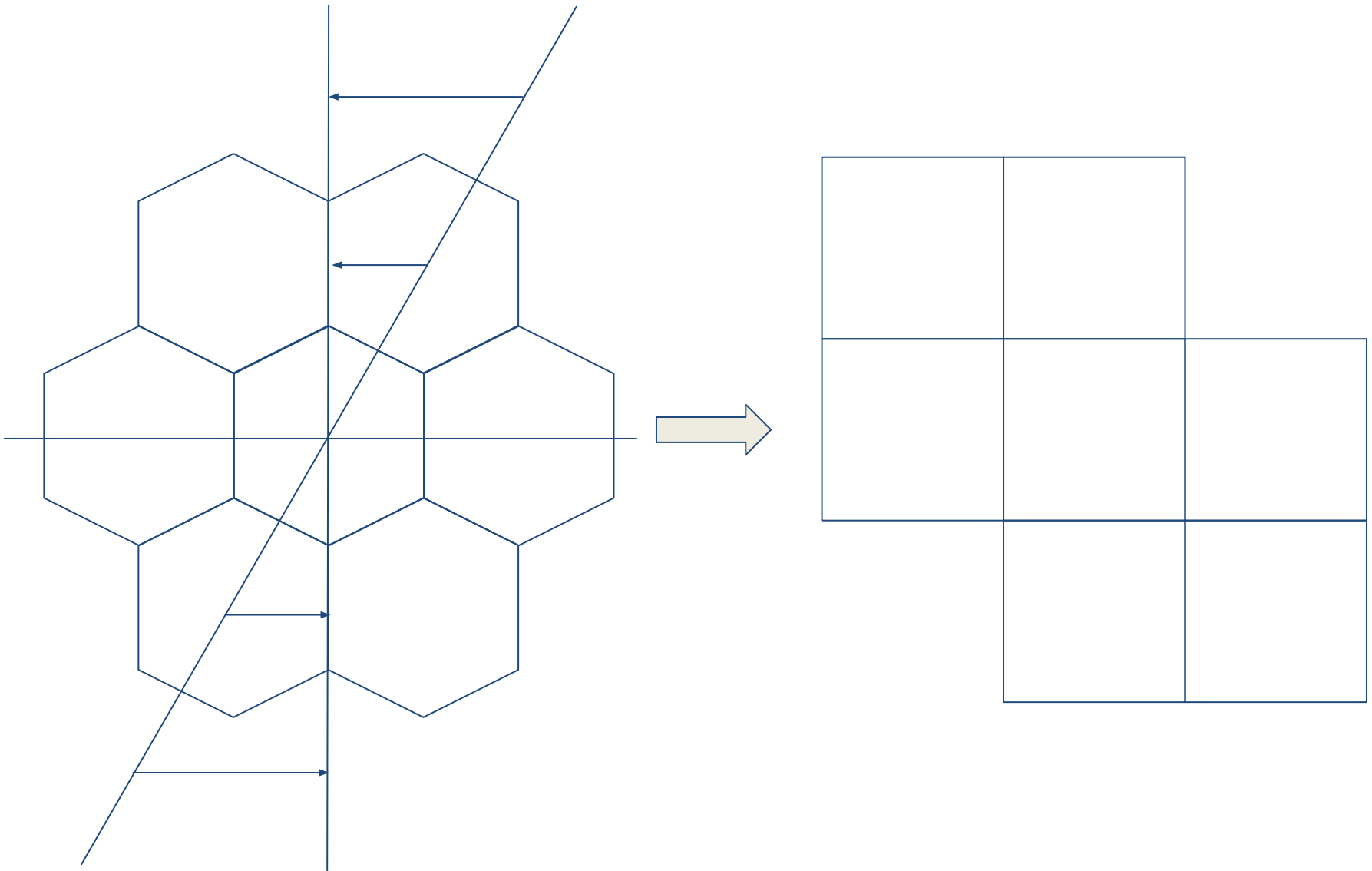
```
def integrate_local_peak(
    waveform: np.ndarray, nb_frames_before_max: int, nb_frames_after_max: int
):
    peaks_idx = np.argmax(waveform, axis=-2, keepdims=True)
    window_idx = peaks_idx + \
        np.arange(-nb_frames_before_max,
                 nb_frames_after_max + 1,
                 dtype=np.int32)[..., np.newaxis]
    waveform_windows = np.take_along_axis(waveform, window_idx, axis=-2)
    charge = waveform_windows.sum(axis=-2)
    return charge
```



Keep pixels that verify 2 conditions:

- condition 1: Pixels are above  $T_1$  and have at least  $N$  neighbors above  $T_1$
- condition 2: Pixels are above  $T_2$  and have at least 1 neighbor passing condition 1

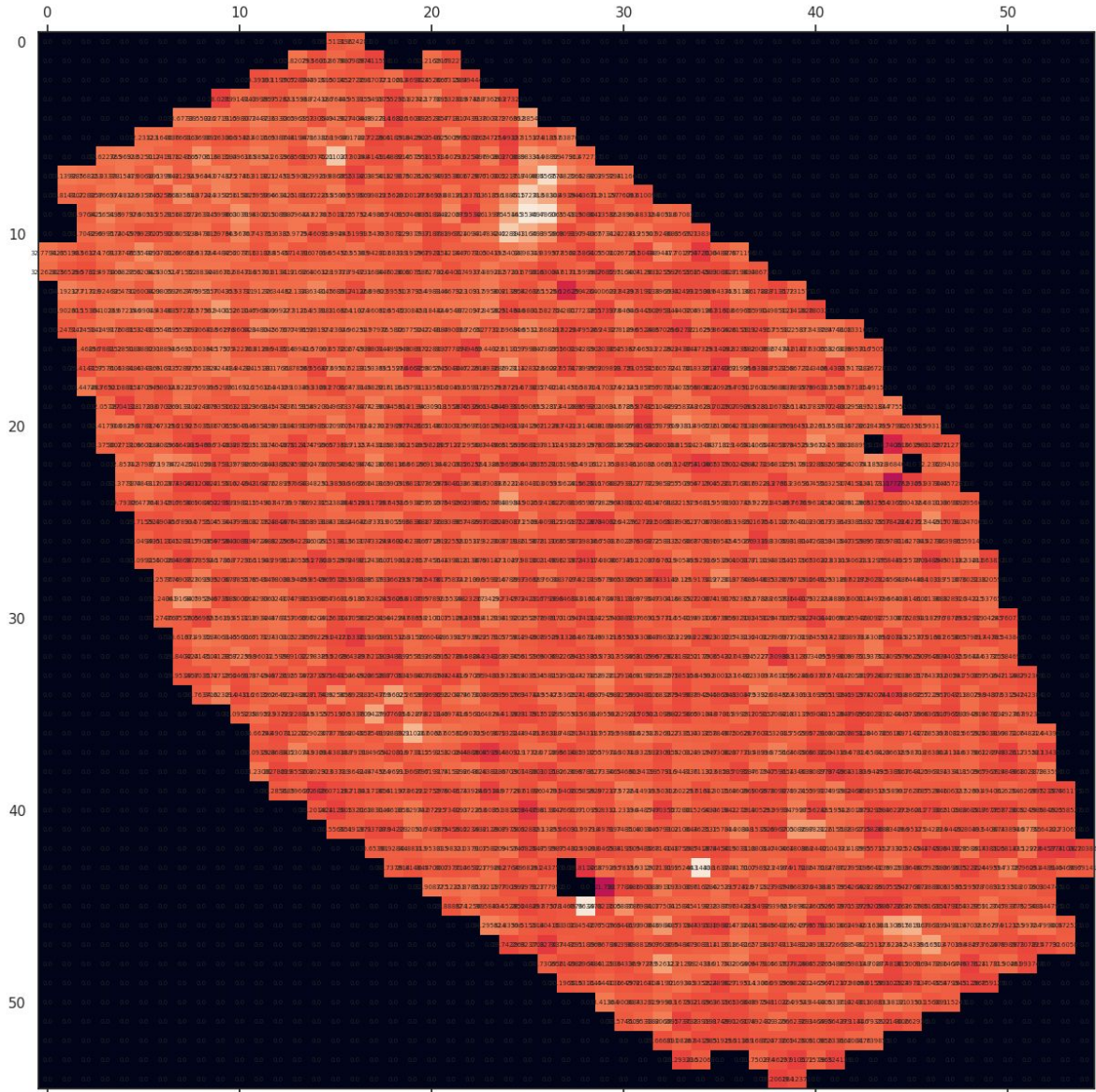






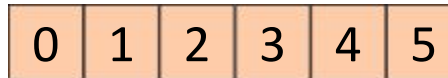
Stencil

1	1	0
1	0	1
0	1	1





```
>>> from numpy.lib.stride_tricks import sliding_window_view
>>> x = np.arange(6)
>>> x.shape
(6,)
>>> v = sliding_window_view(x, 3)
>>> v.shape
(4, 3)
>>> v
array([[0, 1, 2],
       [1, 2, 3],
       [2, 3, 4],
       [3, 4, 5]])
```



```
>>> from numpy.lib.stride_tricks import sliding_window_view
>>> x = np.arange(6)
>>> x.shape
(6,)
>>> v = sliding_window_view(x, 3)
>>> v.shape
(4, 3)
>>> v
array([[0, 1, 2],
       [1, 2, 3],
       [2, 3, 4],
       [3, 4, 5]])
```



```
>>> from numpy.lib.stride_tricks import sliding_window_view
>>> x = np.arange(6)
>>> x.shape
(6,)
>>> v = sliding_window_view(x, 3)
>>> v.shape
(4, 3)
>>> v
array([[0, 1, 2],
       [1, 2, 3],
       [2, 3, 4],
       [3, 4, 5]])
```



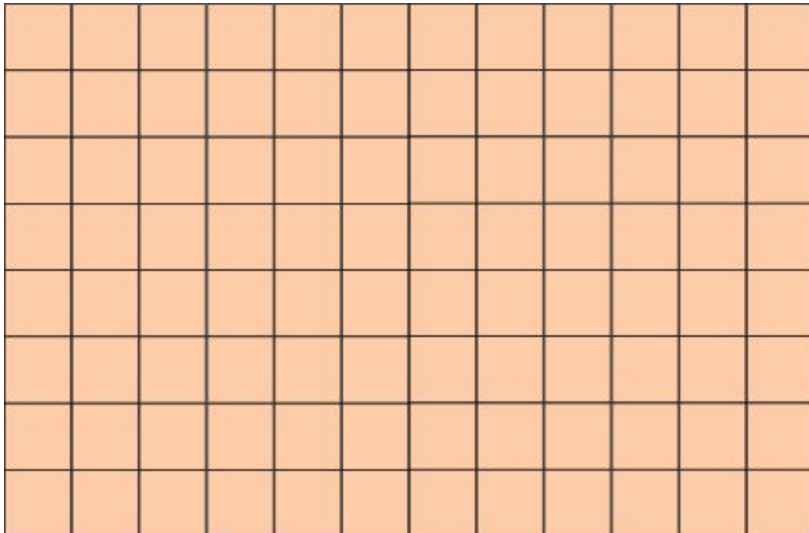
```
>>> from numpy.lib.stride_tricks import sliding_window_view
>>> x = np.arange(6)
>>> x.shape
(6,)
>>> v = sliding_window_view(x, 3)
>>> v.shape
(4, 3)
>>> v
array([[0, 1, 2],
       [1, 2, 3],
       [2, 3, 4],
       [3, 4, 5]])
```



```
>>> from numpy.lib.stride_tricks import sliding_window_view
>>> x = np.arange(6)
>>> x.shape
(6,)
>>> v = sliding_window_view(x, 3)
>>> v.shape
(4, 3)
>>> v
array([[0, 1, 2],
       [1, 2, 3],
       [2, 3, 4],
       [3, 4, 5]])
```

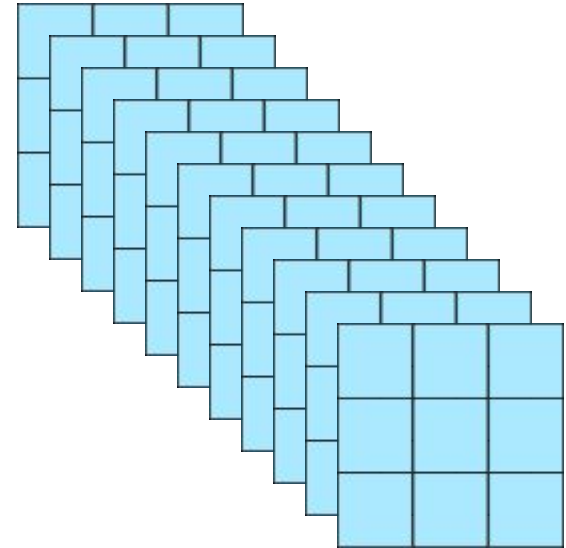
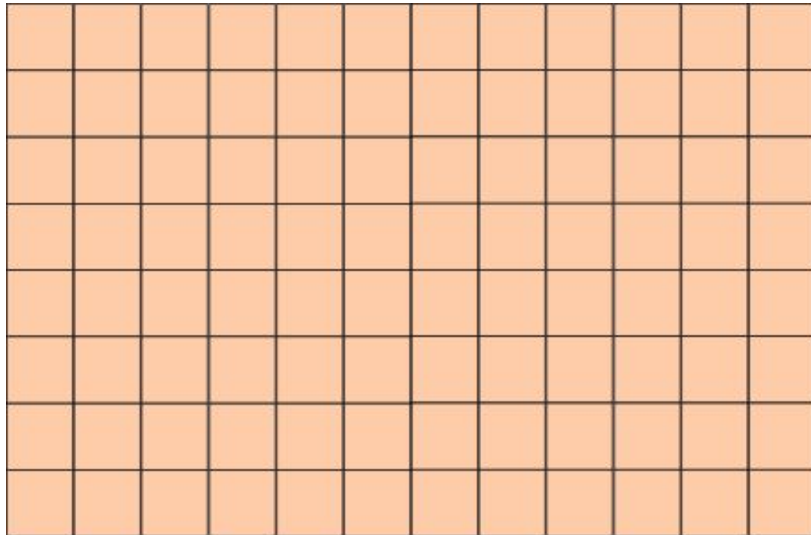


```
neighbors_only_view = sliding_window_view(  
    waveforms_2D_padded,  
    neighbors_stencil.shape  
) * neighbors_stencil  
# condition 1:  
mask = (waveforms_2D >= pixel_threshold) \  
    & ((neighbors_only_view >= pixel_threshold).sum(axis=(-2, -1)) \  
    >= min_number_neighbors)
```

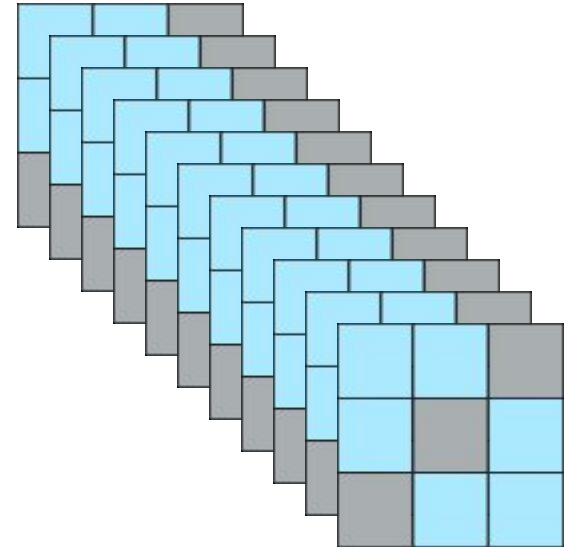
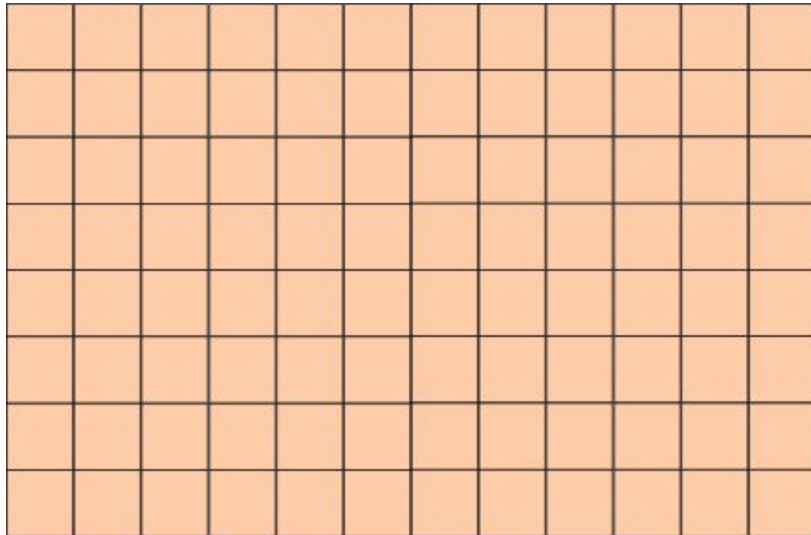




```
neighbors_only_view = sliding_window_view(  
    waveforms_2D_padded,  
    neighbors_stencil.shape  
) * neighbors_stencil  
# condition 1:  
mask = (waveforms_2D >= pixel_threshold) \  
    & ((neighbors_only_view >= pixel_threshold).sum(axis=(-2, -1)) \  
    >= min_number_neighbors)
```



```
neighbors_only_view = sliding_window_view(  
    waveforms_2D_padded,  
    neighbors_stencil.shape  
) * neighbors_stencil  
# condition 1:  
mask = (waveforms_2D >= pixel_threshold) \  
    & ((neighbors_only_view >= pixel_threshold).sum(axis=(-2, -1)) \  
    >= min_number_neighbors)
```

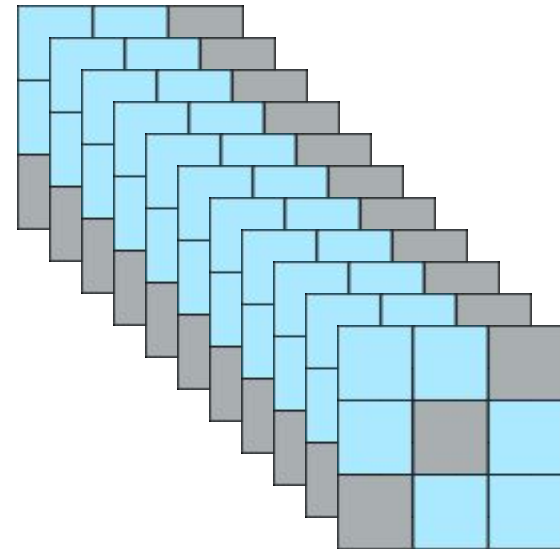
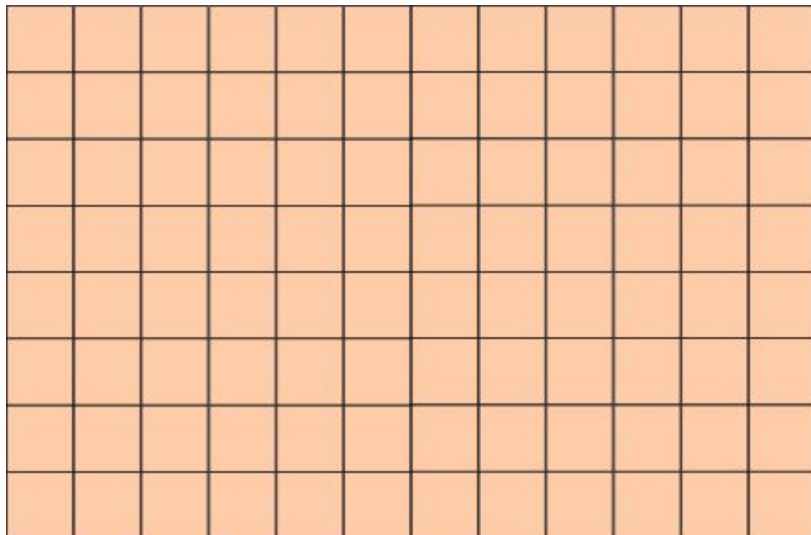


```

neighbors_only_view = sliding_window_view (
    waveforms_2D_padded,
    neighbors_stencil.shape
) * neighbors_stencil
# condition 1:
mask = (waveforms_2D >= pixel_threshold) \
    & ((neighbors_only_view >= pixel_threshold).sum(axis=(-2, -1)) \
    >= min_number_neighbors)
    
```

Pixel above threshold

N neighbors above threshold

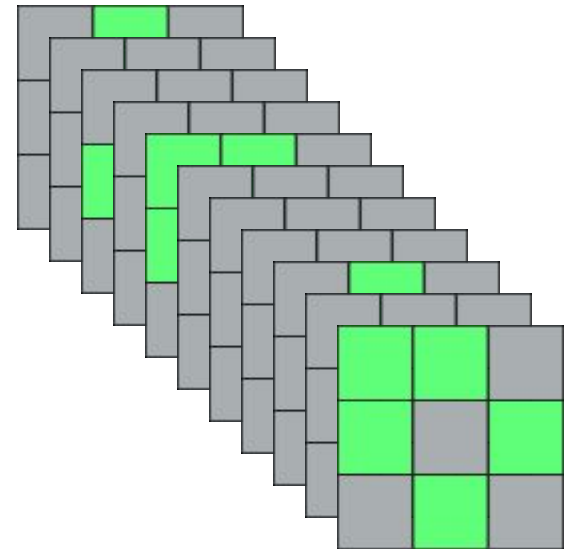
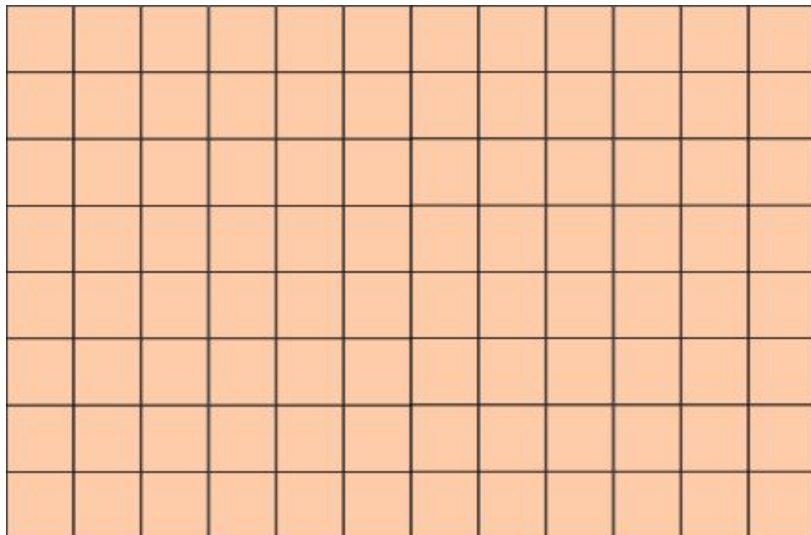


```

neighbors_only_view = sliding_window_view (
    waveforms_2D_padded,
    neighbors_stencil.shape
) * neighbors_stencil
# condition 1:
mask = (waveforms_2D >= pixel_threshold) \
    & ((neighbors_only_view >= pixel_threshold).sum(axis=(-2, -1)) \
    >= min_number_neighbors)
    
```

Pixel above threshold

N neighbors above threshold

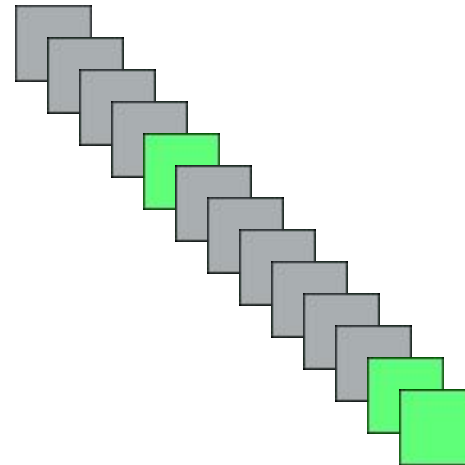
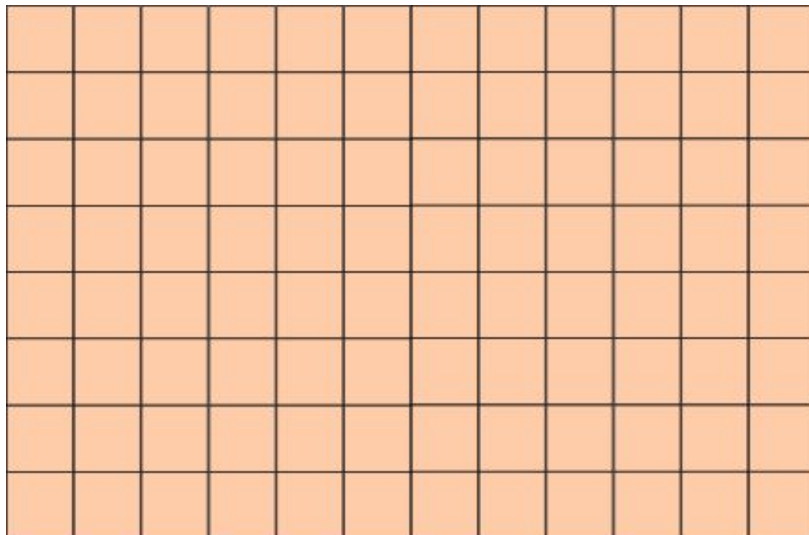


```

neighbors_only_view = sliding_window_view (
    waveforms_2D_padded,
    neighbors_stencil.shape
) * neighbors_stencil
# condition 1:
mask = (waveforms_2D >= pixel_threshold) \
    & ((neighbors_only_view >= pixel_threshold).sum(axis=(-2, -1)) \
    >= min_number_neighbors)
    
```

Pixel above threshold

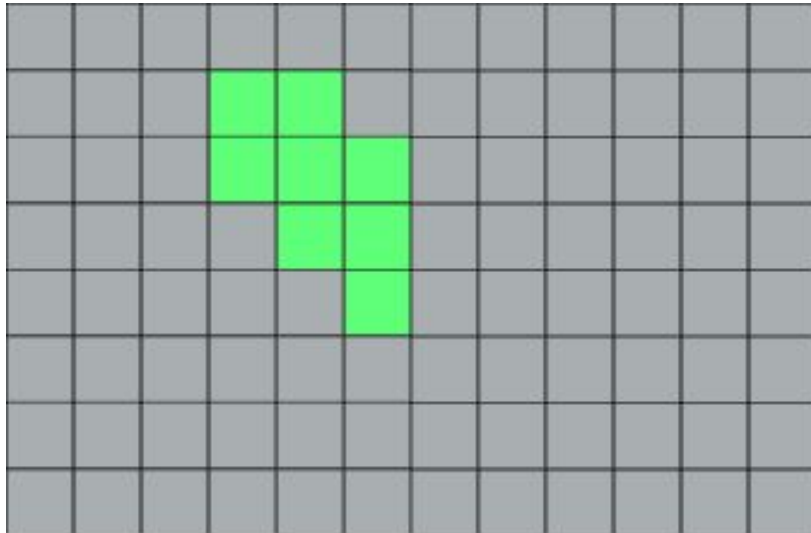
N neighbors above threshold



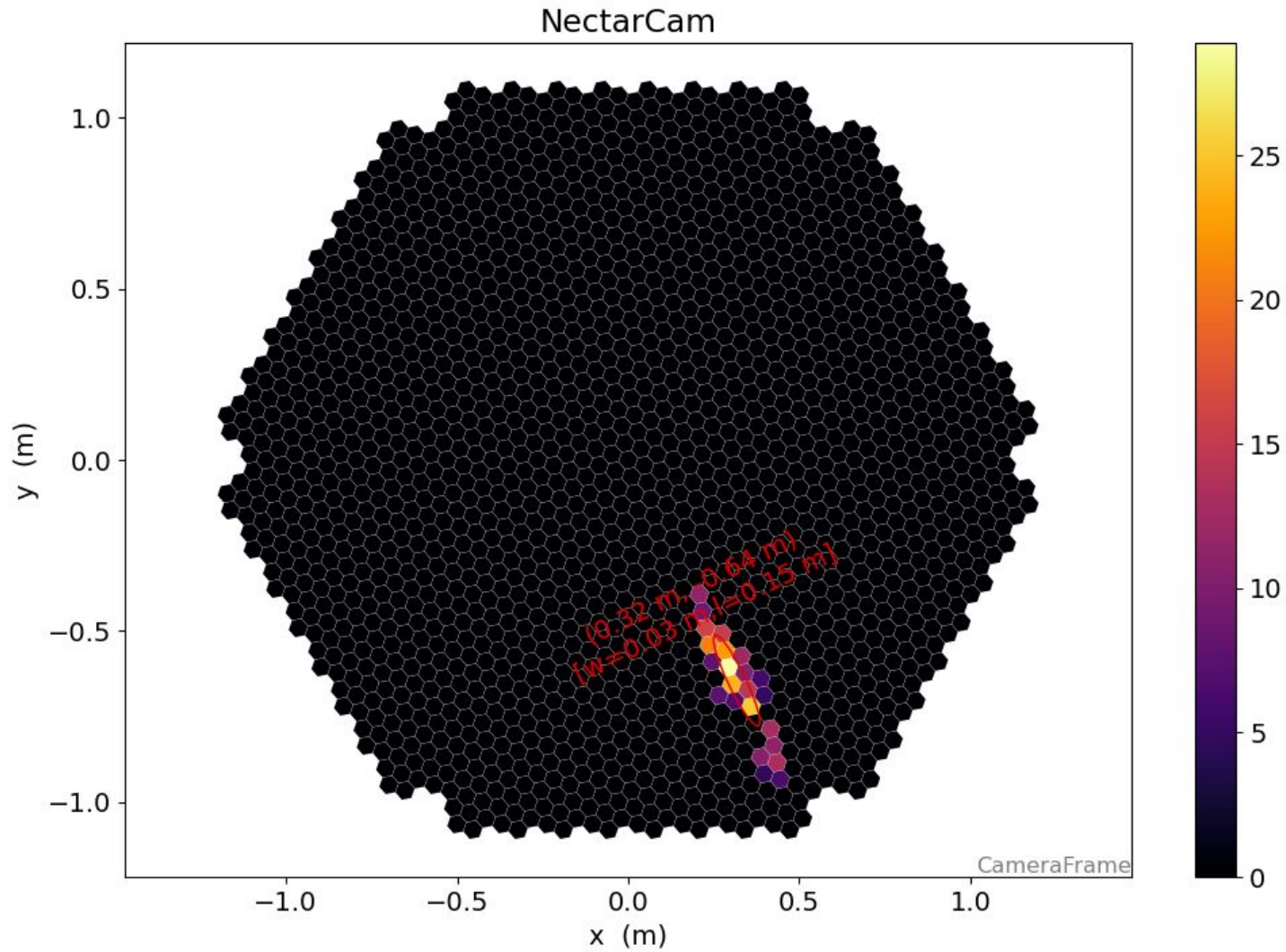
```
neighbors_only_view = sliding_window_view(  
    waveforms_2D_padded,  
    neighbors_stencil.shape  
) * neighbors_stencil  
# condition 1:  
mask = (waveforms_2D >= pixel_threshold) \  
    & ((neighbors_only_view >= pixel_threshold).sum(axis=(-2, -1)) \  
    >= min_number_neighbors)
```

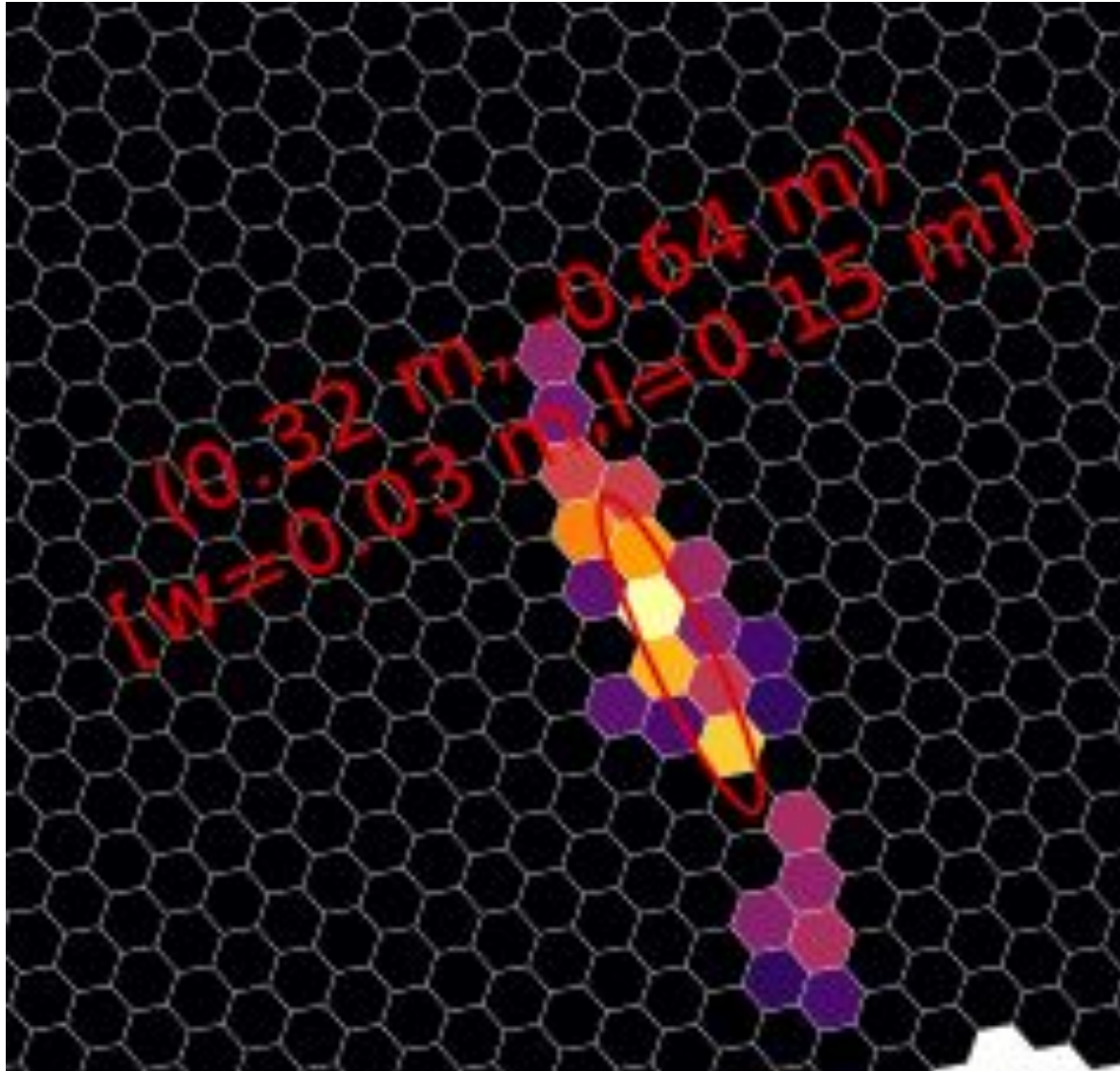
Pixel above threshold

N neighbors above threshold









Ellipse fit: Principal Component Analysis transformation

- Estimation of the covariance matrix of the image pixels

```
def hillas(pix_x, pix_y, image_batch):  
    for image in image_batch:  
        # center of gravity  
        cog_x = np.average(pix_x, weights=image)  
        cog_y = np.average(pix_y, weights=image)  
        # PCA  
        delta_x = pix_x - cog_x  
        delta_y = pix_y - cog_y  
        cov = np.cov(delta_x, delta_y, aweights=image, ddof=0)  
        eig_vals, eig_vecs = np.linalg.eigh(cov)  
        # width and length are eigen values of the PCA  
        width, length = np.sqrt(eig_vals)
```

← sequential operations

Issue: Numpy covariance, eigenvalue functions do not handle batches

$$\text{Var}(\vec{X}) = \begin{pmatrix} \text{Var}(X_1) & \text{Cov}(X_1, X_2) & \cdots & \text{Cov}(X_1, X_p) \\ \text{Cov}(X_2, X_1) & \ddots & \cdots & \vdots \\ \vdots & \vdots & \ddots & \vdots \\ \text{Cov}(X_p, X_1) & \cdots & \cdots & \text{Var}(X_p) \end{pmatrix}$$

$$\mathbb{E}[X] = \sum_{i=1}^{\infty} x_i p_i,$$

$$\text{Var}(X) = \mathbb{E}[X^2] - \mathbb{E}[X]^2$$

$$\text{cov}(X, Y) = \mathbb{E}[(X - \mathbb{E}[X])(Y - \mathbb{E}[Y])]$$

```
def hillas(images, position):
    intensity = waveform.sum(axis=-1)

    moments = (
        waveform[..., np.newaxis, :] # shape (N_batch, 1, N_pixels) to broadcast for all moments
        * np.stack([position[0], position[1], position[0] ** 2, position[1] ** 2, position.prod(axis=0)])
    ).sum(axis=-1) / intensity[..., np.newaxis]

    E_x_y_square = moments[..., 0:2] ** 2
    variance = moments[..., 2:4] - E_x_y_square
    covariance = moments[..., 4] - moments[..., 0:2].prod(axis=-1)
    varx_plus_vary = variance.sum(axis=-1)

    char_pol_delta = np.sqrt(varx_plus_vary**2 + 4.0 * (covariance**2 - variance.prod(axis=-1)))
    eigenValueHigh = (varx_plus_vary + char_pol_delta) / 2.0
    eigenValueLow = (varx_plus_vary - char_pol_delta) / 2.0

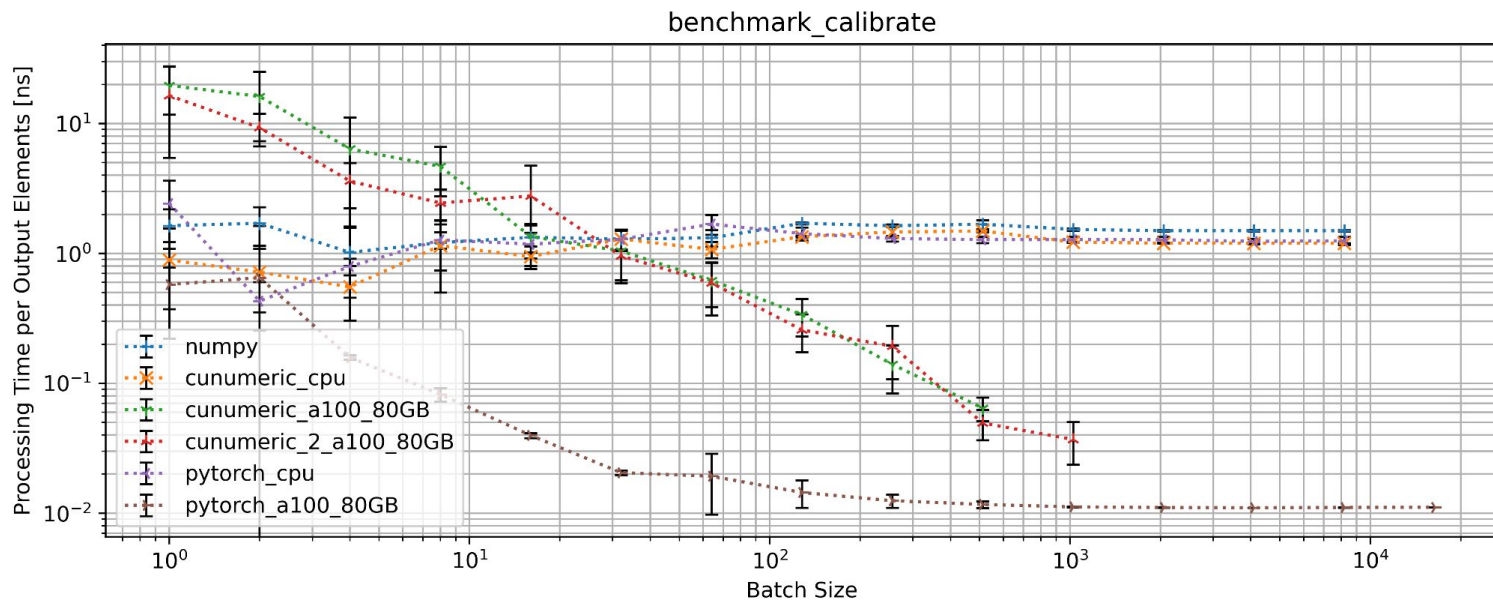
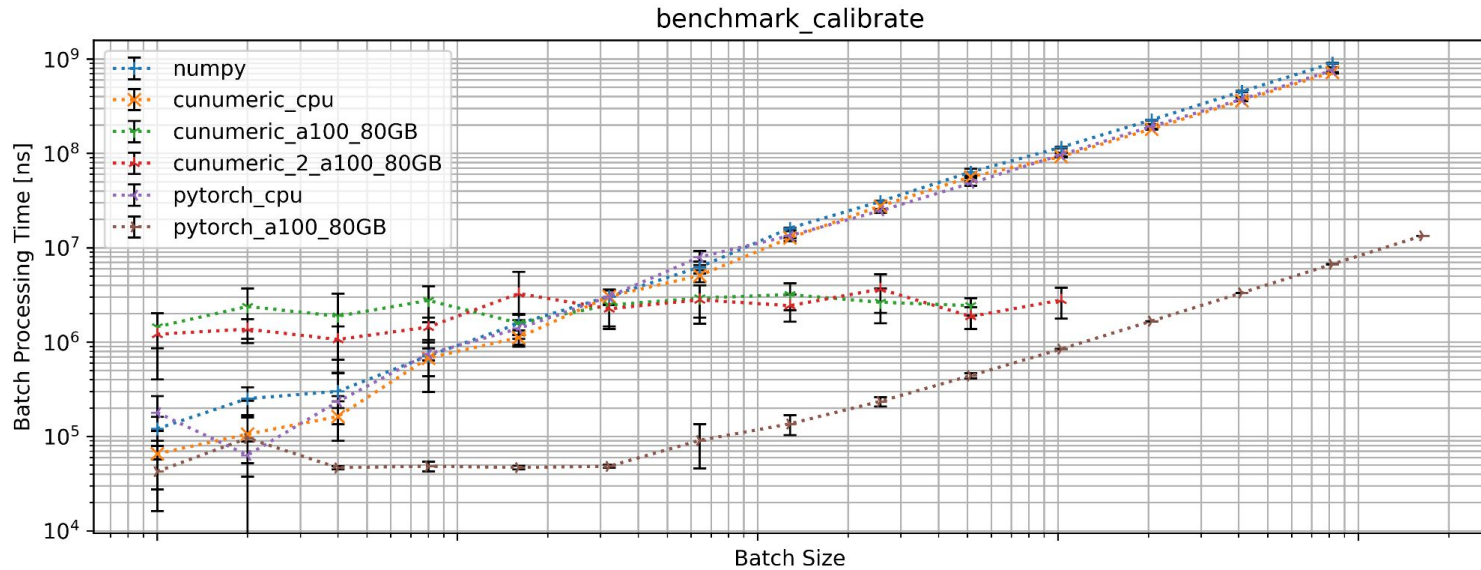
    length = np.sqrt(eigenValueHigh)
    width = np.sqrt(eigenValueLow)
```

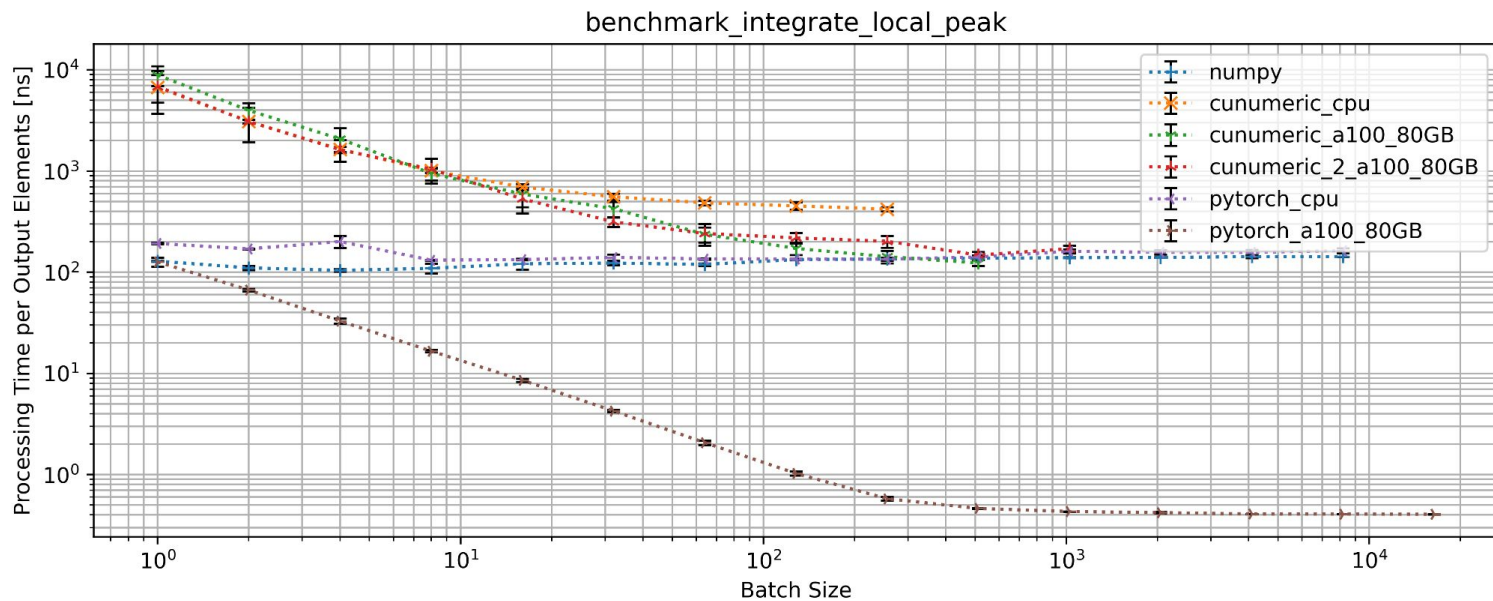
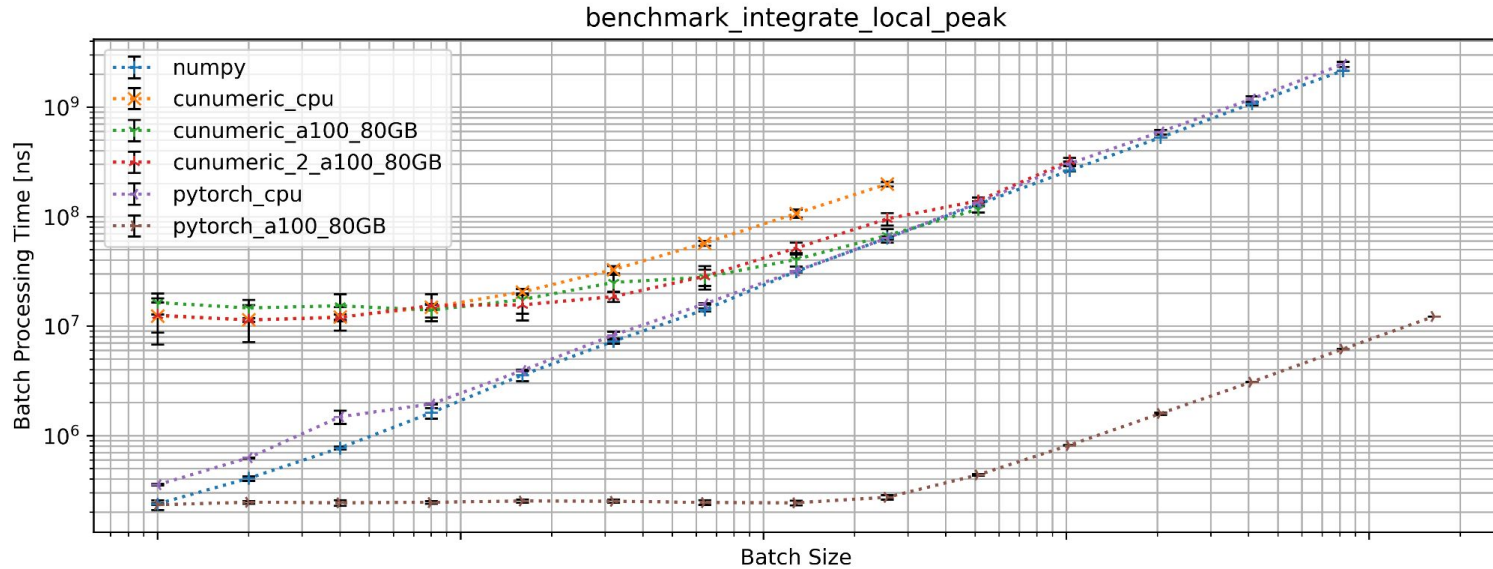
- Installation
  - version 24.03 from source (first legate.core, then cuNumeric)
  - dependencies installed with cuda (cuTensor cuTensor-cuda missing)
  - requires machine with gpu to enable gpu support
- Translation
  - `import cunumeric as np`
  - Differences in API of `cunumeric.random.default_rng`
  - no “stride\_tricks” -> no cleaning
- Execution
  - using legate driver with several cpus, gpus on same node
  - much more memory used
  - crash when not enough memory is available (even on cpu)

“Whenever possible, use the `out` parameter in the APIs, to avoid allocating an intermediate array in our implementation.” cuNumeric guidelines

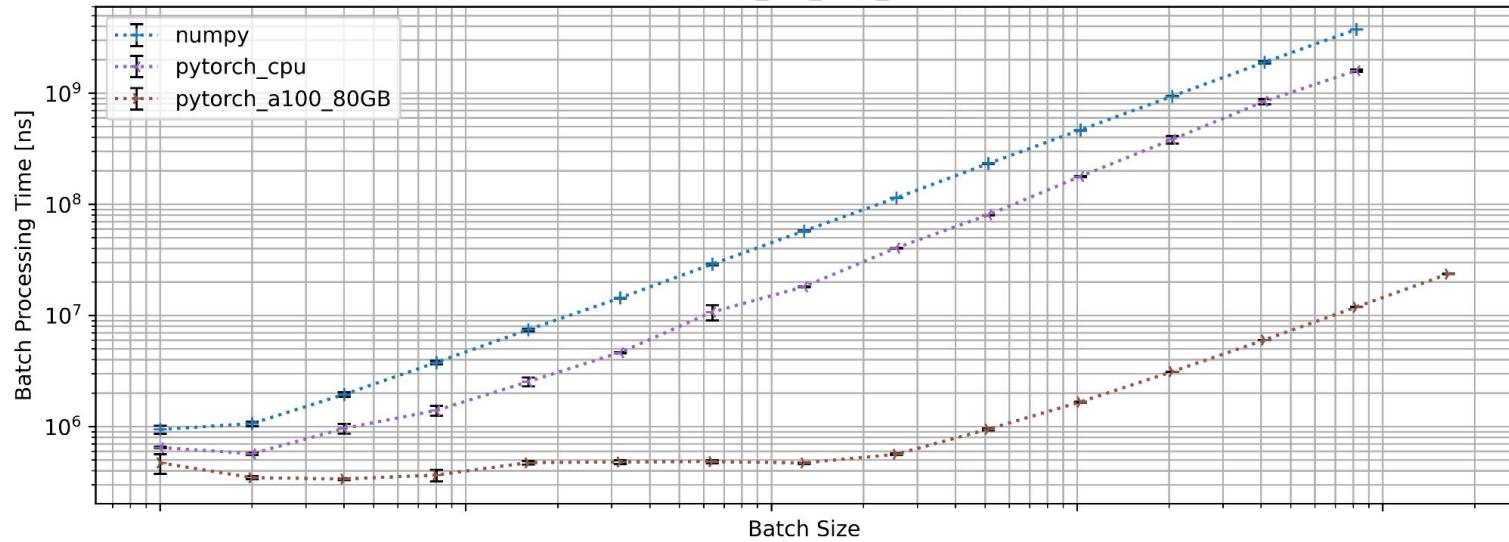
- Installation
  - from pytorch conda channel, with cuda 12.1
- Translation
  - tedious but easy as API is very similar with some name changes
    - `np.*` -> `torch.*`
    - `np.ndarray` -> `torch.Tensor`
    - `np.array(...)` -> `torch.tensor(...)`
    - slightly different padding
    - no `striding_window_view` but `as_strided` is supported
    - introduce `torch.device`
    - fix `dtype` to `float32` or `int32`
- Execution
  - multi-cpu (mono node)
  - mono gpu



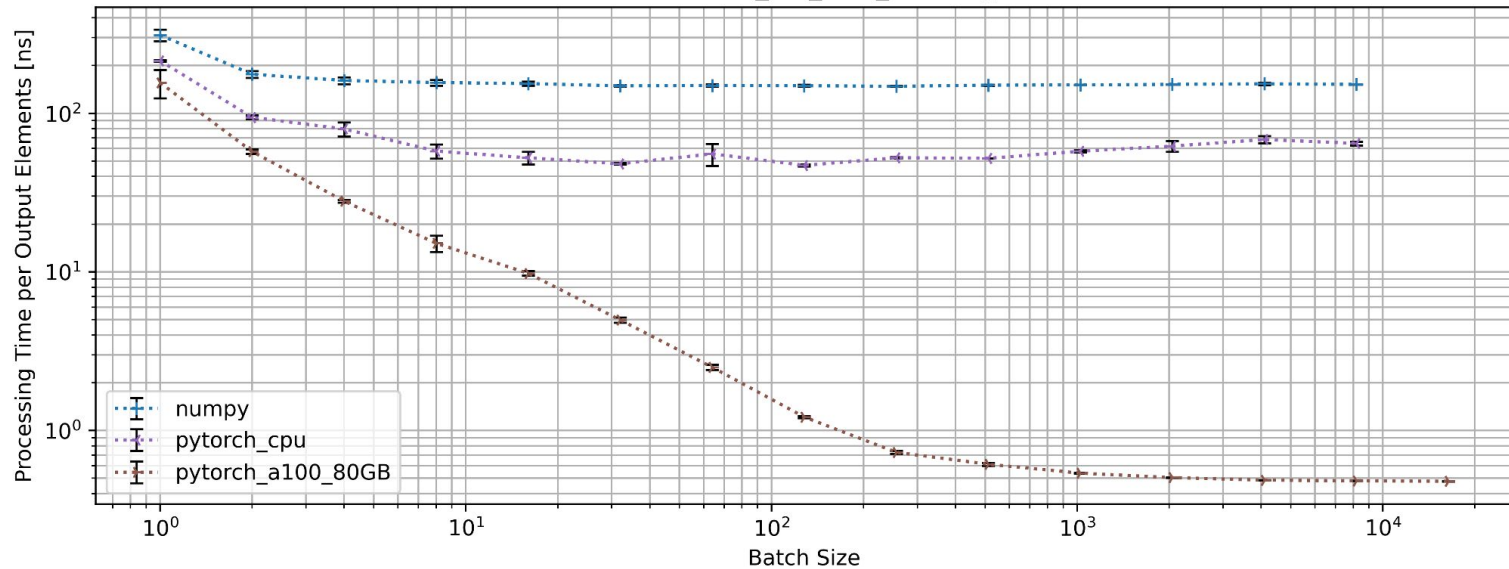


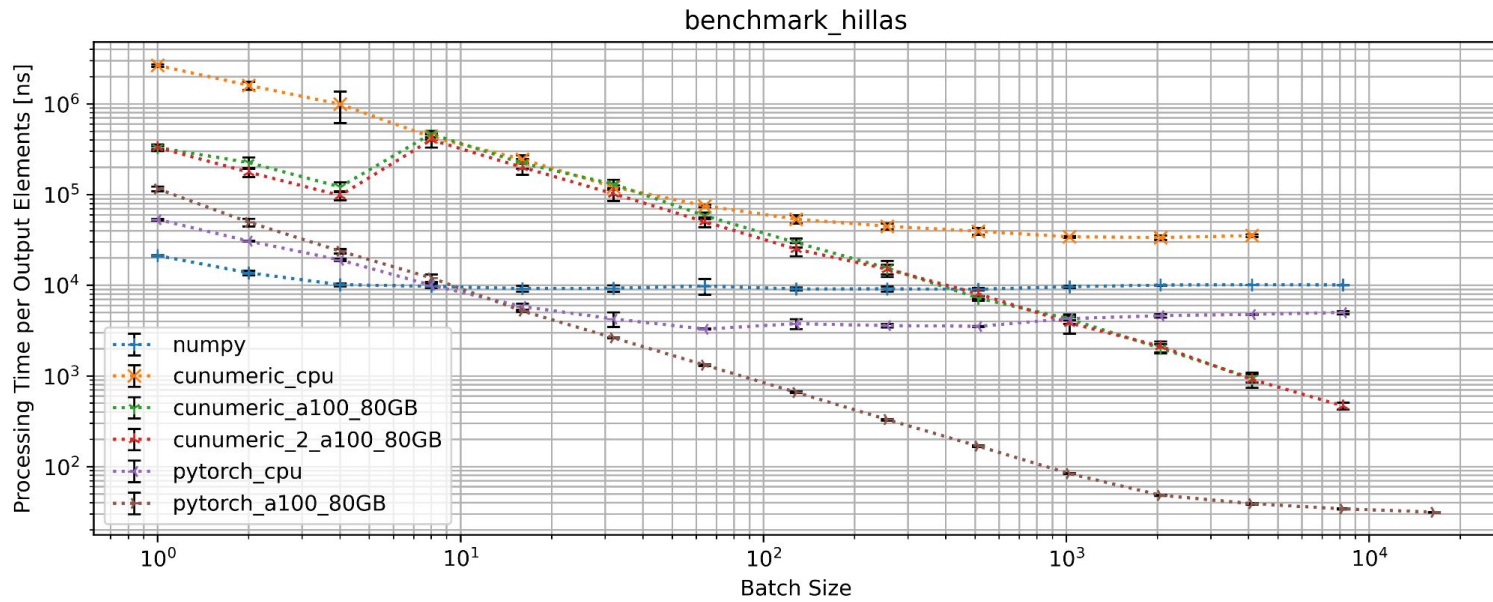
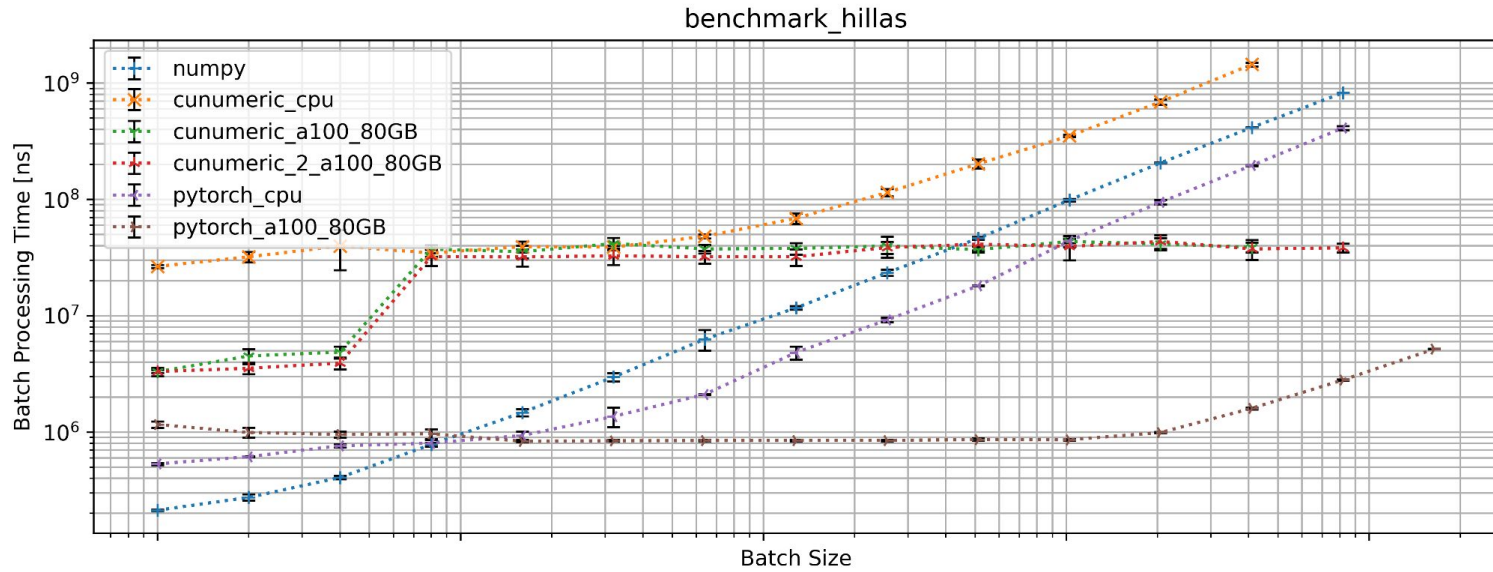


benchmark\_tail\_cuts\_cleaning



benchmark\_tail\_cuts\_cleaning





- Possible to have a single code base with *decent* performance on cpu and gpu implemented with python
- cuNumeric
  - still in development
  - doesn't deliver all its promises (yet?)
  - transparent scaling to multi-gpu (multi-nodes?)
- pytorch
  - mature library
  - API very similar to numpy
  - good performances on cpu and gpu

#### Next steps:

- pre-allocated memory layout
- benchmark cpu to gpu data transfer
  - cuNumeric GPU direct storage
  - pytorch: mitigate transfer cost behind more computation ?

<https://gitlab.in2p3.fr/CTA-LAPP/rta/pyhiperta/>