

Profile, Optimize, Repeat

One Core Is All You Need™

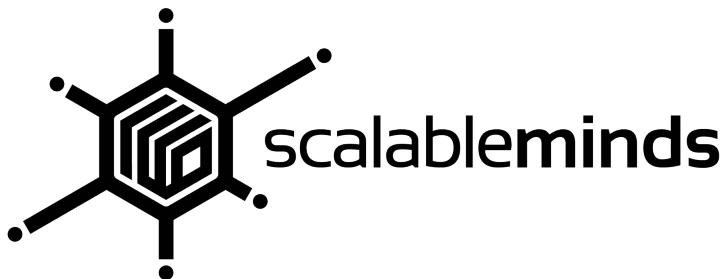
Jonathan Striebel & Valentin Nieper



EuroPython2024

Hi, we're

Valentin Nieper & Jonathan Striebel



agnostics



@jostriebel



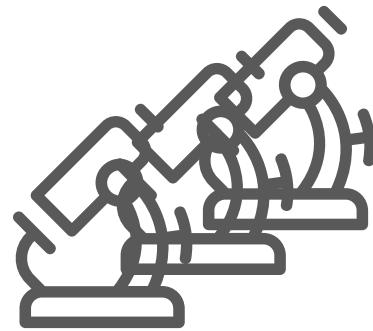
codical.org

Data Science Workloads Grow

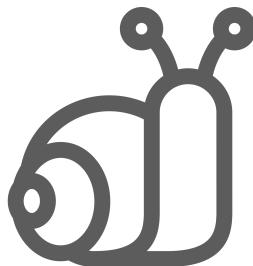


Analyzing Biomedical Data

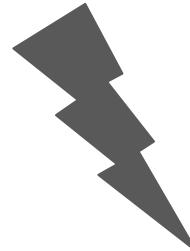
Data Science Workloads Grow



Analyzing Biomedical Data
even^more

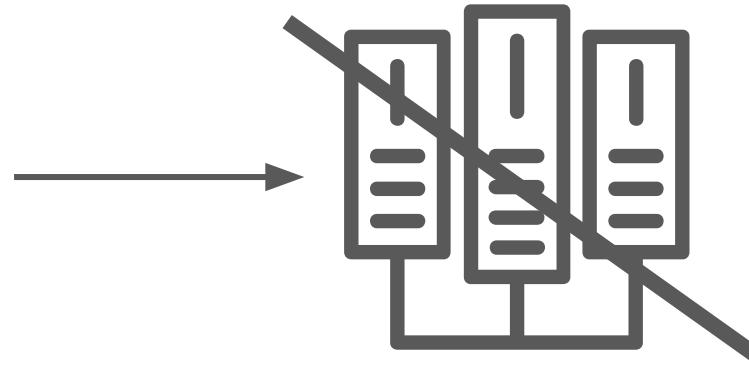


Too slow



OOM :-(

Data Science Workloads Grow



Don't scale up/out (yet)

Data Science Workloads Grow

- Single core improvements pay off, also when parallelized later
- Parallelization needs resources
(cores, memory, cluster-nodes, money, time)
- Code may be hard / impossible to parallelize
(e.g. when reducing results from a map-reduce)

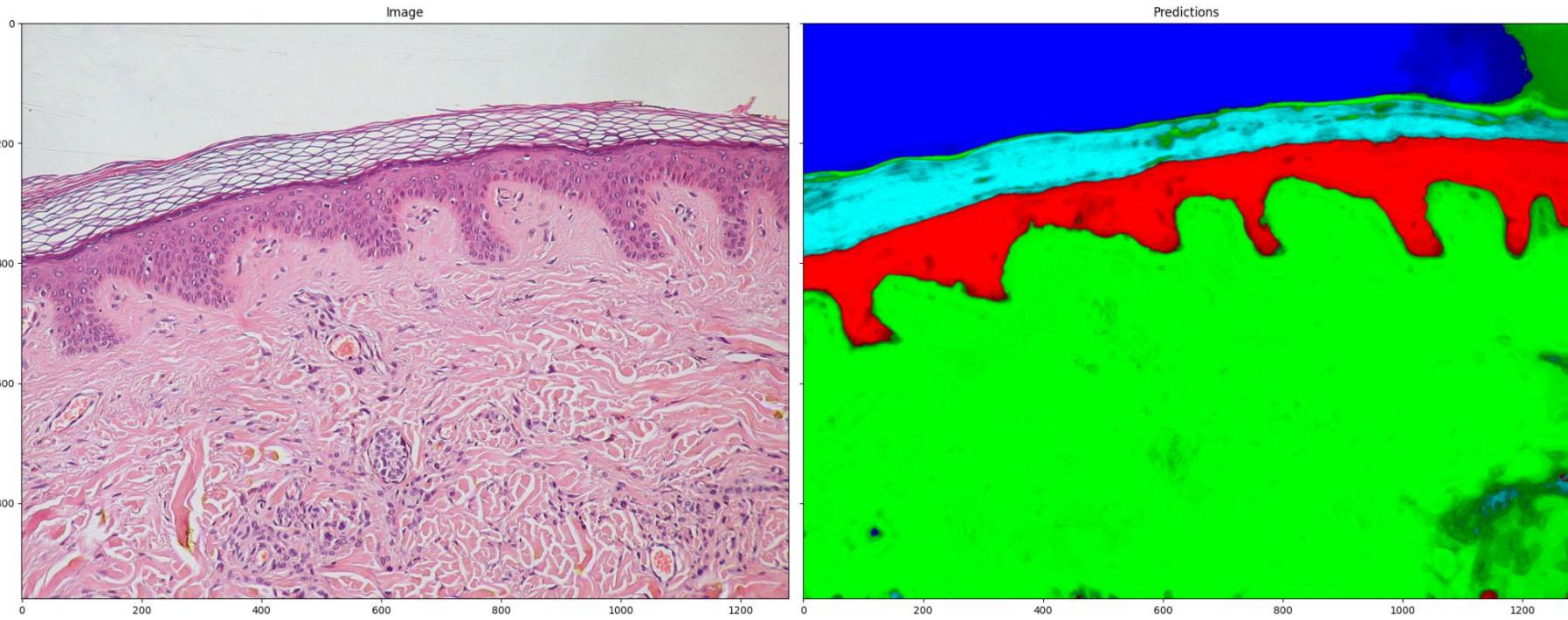
Don't scale up/out (yet)

One Core Is All You Need™

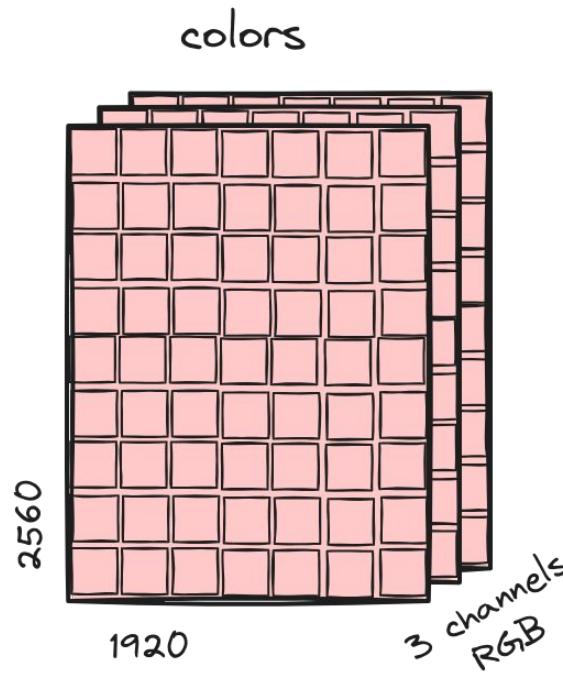
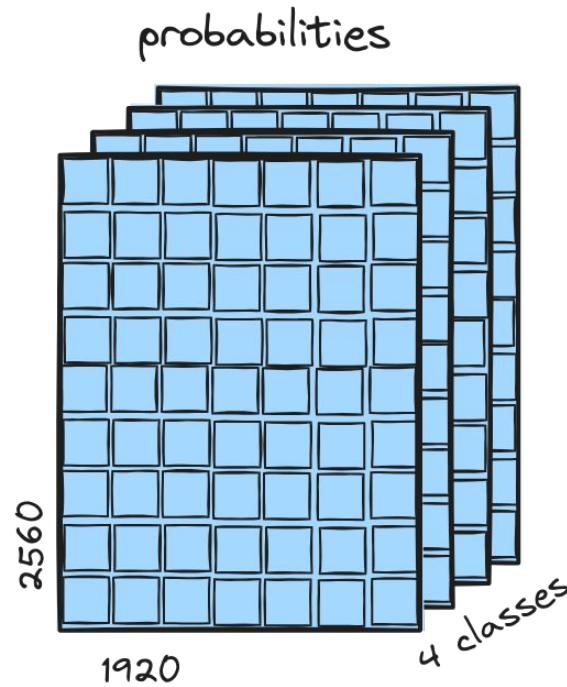
Example: Colorcode Predictions

Based on scikit-image tutorial

[“Trainable segmentation using local features and random forests”](#)



Example: Colorcode Predictions



Example: Colorcode Predictions

```
def load_probabilities(folder):
    arrays = [np.loadtxt(folder / f"{c}.csv") for c in range(4)]
    return np.stack(arrays, axis=-1)
```

Example: Colorcode Predictions

```
def colorcode_probabilities(probabilities):
    class_colors = [[0, 0, 255], [0, 255, 0], [255, 0, 0], [0, 255, 255]]
    colored_probabilities = []
    for probabilities_row in probabilities:
        colored_probabilities_row = []
        for class_probabilities in probabilities_row:
            class_index = None
            max_prob = 0
            for i, prob in enumerate(class_probabilities):
                if prob > max_prob:
                    max_prob = prob
                    class_index = i
            colored_probabilities_row.append([c * max_prob for c in class_colors[class_index]])
        colored_probabilities.append(colored_probabilities_row)
    return colored_probabilities
```

Example: Colorcode Predictions

Profile, Optimize, Repeat

15s → 150ms

800MB → 80MB



Profile

- Measure
- Speed & Memory Bottlenecks
- to Identify & Analyse them
- for Mitigation

What?

Why?



Time/CPU Profiling

Slow Code

- Where?
- How slow?
- Why?



Memory Profiling

High Memory Usage

- Where?
- How much?
- Why?



Instrumenting vs Sampling

Deterministic Instrumenting

- Measure **each function or code block** of a program
- Potential overhead
- Inaccuracies for calls with lots of instrumentation inside

Statistical Sampling

- **Periodically sample** the program's state
- Can miss brief invocations
(usually fine for timing, problematic for short memory spikes)



Scalene

CPU, Memory, GPU profiling

sampling-based

nice visualization



More

Profiler	Link	Time	Mem	Viz	Type	Other
cProfile	🔗		X	X	instrumenting	Viz via SnakeViz 🔗 or Tuna 🔗
py-spy	🔗		X		sampling	
memray	🔗	X			instrumenting	
Scalene	🔗				sampling	also GPU time
Austin	🔗				sampling	can also be used with pprof 🔗
pyinstrument	🔗		X		sampling	
line_profiler	🔗		X	X	instrumenting	
Fil	🔗	X			instrumenting	
Guppy3	🔗	X		X	instrumenting	
Yappi	🔗		X	X	instrumenting	
Palanteer	🔗				instrumenting	also supports C++
vprof	🔗				sampling	
VizTracer	🔗		X		instrumenting	
sciagraph	🔗				sampling	proprietary profiler



codical.org

Run scalene

main.py

```
from scalene import scalene_profiler  
scalene_profiler.start()  
CODE  
scalene_profiler.stop()
```

```
scalene --off --cli --cpu main.py --- v0 --data data/m
```

15s runtime

Line	Time	Python	native	system	
1					
2					
3					
4					
5	24%				
6					
7					
8					
9					
10					
11					
12					
13					
14					
15					
16					
17	6%				
18					
19					
20					
21	41%		9%	1%	
22					
23					
24					
25					

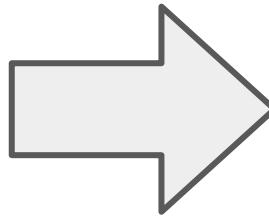
Efficient IO

Text-based

csv

json

yaml



Binary format

hdf5

npy

parquet

pickle

sqlite

zarr

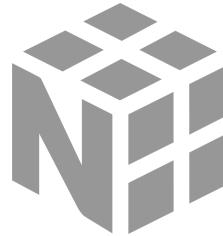
...

see also

github.com/mverleg/array_storage_benchmark

Line	Time	-----	-----		10.6s runtime
	Python	native	system		
1				<pre>import h5py</pre>	
2					
3				<pre>def load_probabilities(folder):</pre>	
4				<pre> with h5py.File(folder / "data.hdf5", "r") as f:</pre>	
5				<pre> return f["probabilities"][:]</pre>	
6					
7				<pre>def colorcode_probabilities(probabilities):</pre>	
8				<pre> class_colors = [[0, 0, 255], [0, 255, 0], [255, 0, 0], [0, 255, 255]]</pre>	
9				<pre> colored_probabilities = []</pre>	
10				<pre> for probabilities_row in probabilities:</pre>	
11				<pre> colored_probabilities_row = []</pre>	
12				<pre> for class_probabilities in probabilities_row:</pre>	
13				<pre> class_index = None</pre>	
14				<pre> max_prob = 0</pre>	
15	9%			<pre> for i, prob in enumerate(class_probabilities):</pre>	
16				<pre> if prob > max_prob:</pre>	
17				<pre> max_prob = prob</pre>	
18				<pre> class_index = i</pre>	
19	52%	13%	3%	<pre> colored_probability = [c * max_prob for c in class_colors[class_index]]</pre>	
20				<pre> colored_probabilities_row.append(colored_probability)</pre>	
21				<pre> colored_probabilities.append(colored_probabilities_row)</pre>	
22				<pre> return colored_probabilities</pre>	

Vectorization



NumPy

- compact in-memory representation for homogeneous data
 - less powerful than Python types, but less overhead
 - optimized primitives for
 - add, multiply, ...
 - argmax, unique, ...
 - fancy indexing
- NumPy core is written in C

Vectorization

```
def colorcode_probabilities(probabilities):
    predicted_classes = np.argmax(probabilities, axis=2)
    max_prob = np.max(probabilities, axis=2)
    class_colors = np.array([[[0, 0, 255], [0, 255, 0], [255, 0, 0], [0, 255, 255]]], dtype=np.uint8)
    colored_probabilities = class_colors[predicted_classes]
    colored_probabilities = colored_probabilities * max_prob[..., np.newaxis]
    return colored_probabilities.astype(np.uint8)
```

Line	Time	Python	native	system
1				
2				
3				
4				
5				
6				
7				
8				
9	1%	4%	1%	
10		26%	16%	
11				
12				
13				
14		17%	32%	

```
import numpy as np
import h5py

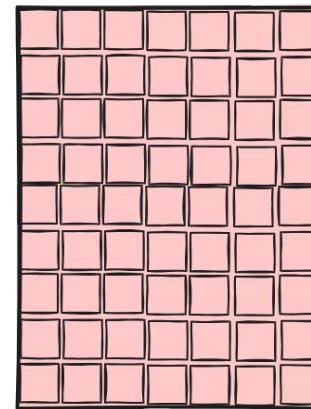
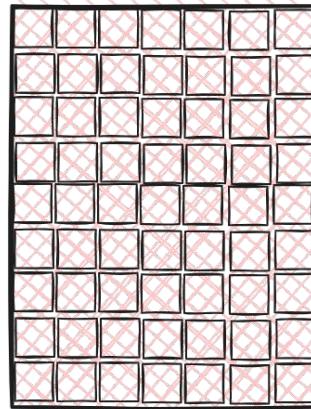
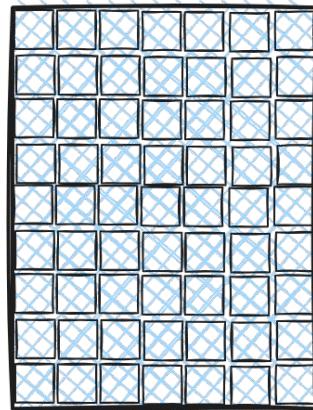
def load_probabilities(folder):
    with h5py.File(folder / "data.hdf5", "r") as f:
        return f["probabilities"][:]

def colorcode_probabilities(probabilities):
    predicted_classes = np.argmax(probabilities, axis=2)
    max_prob = np.max(probabilities, axis=2)
    class_colors = np.array([[0, 0, 255], [0, 255, 0], [255, 0, 0], [0, 255, 255]]), dtype=np.uint8
    colored_probabilities = class_colors[predicted_classes]
    colored_probabilities = colored_probabilities * max_prob[..., np.newaxis]
    return colored_probabilities.astype(np.uint8)
```

Vectorization

```
def colorcode_probabilities(probabilities):
    predicted_classes = np.argmax(probabilities, axis=2)
    max_prob = np.max(probabilities, axis=2)
    class_colors = np.array([[[0, 0, 255], [0, 255, 0], [255, 0, 0], [0, 255, 255]]], dtype=np.uint8)
    colored_probabilities = class_colors[predicted_classes]
    colored_probabilities = colored_probabilities * max_prob[..., np.newaxis]
    return colored_probabilities.astype(np.uint8)
```

157MB



15MB

352MB peak mem

Memory Python	peak	timeline/%	Copy (MB/s)	
				<pre>import numpy as np import h5py # Change: all numpy def load_probabilities(folder): with h5py.File(folder / "data.hdf5", "r") as f: return f["probabilities"][:] def colordcode_probabilities(probabilities): predicted_classes = np.argmax(probabilities, axis=2) max_prob = np.max(probabilities, axis=2) class_colors = np.array([[0, 0, 255], [0, 255, 0], [255, 0, 0], [0, 255, 255]]), dtype=np.uint8) colored_probabilities = class_colors[predicted_classes] colored_probabilities = colored_probabilities * max_prob[..., np.newaxis] return colored_probabilities.astype(np.uint8)</pre>
151M	151M	41%		
68%	37M	10%		
36%	37M	10%		
7%	14M	4%	80	
	113M	31%		
	14M	4%		

Reduce Memory Usage

- Optimize data types

uint8	0 to 255	1 byte
int32	-2147483648 to 2147483647	4 bytes
float32	-3.4e+38 to 3.4e+38	4 bytes
float64	-1.79e+308 to 1.79e+308	8 bytes

- Manual GC
- Manual loop-unrolling (avoid broadcasting)

Optimized Numpy

```
def colorcode_probabilities(probabilities):
    # use uint8 to reference the classes (we only have 4)
    predicted_classes = np.argmax(probabilities, axis=2).astype(np.uint8)
    # Do no recompute max, use argmax result to index
    max_prob = np.take_along_axis(probabilities, predicted_classes[..., None], axis=2).squeeze()
    # manual GC, we don't need the probabilities anymore
    del probabilities
    # use float32 to save memory
    max_prob = max_prob.astype(np.float32)
    class_colors = np.array([object: [[0, 0, 255], [0, 255, 0], [255, 0, 0], [0, 255, 255]], dtype=np.uint8])
    colored_probabilities = class_colors[predicted_classes]
    # avoid broadcasting by unrolling the RGB loop
    colored_probabilities[:, :, 0] = colored_probabilities[:, :, 0] * max_prob
    colored_probabilities[:, :, 1] = colored_probabilities[:, :, 1] * max_prob
    colored_probabilities[:, :, 2] = colored_probabilities[:, :, 2] * max_prob
    return colored_probabilities
```

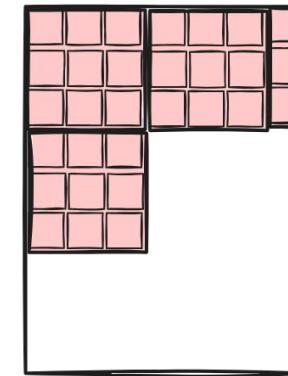
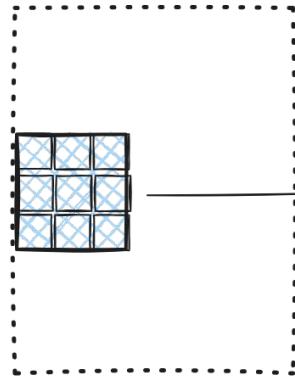
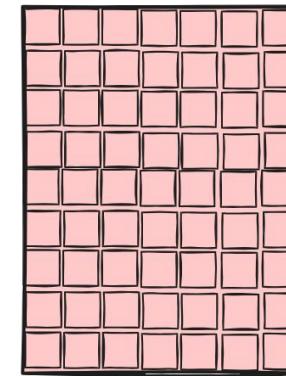
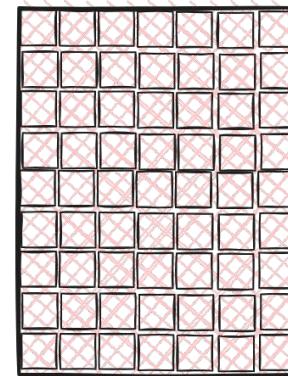
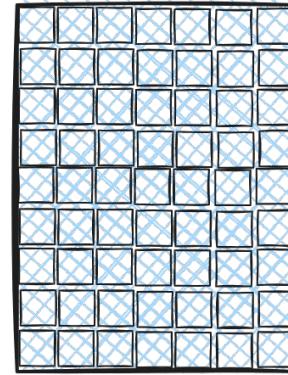
211MB peak mem

Memory Python	peak	timeline/%	Copy (MB/s)	
				<pre>import numpy as np import h5py def load_probabilities(folder): with h5py.File(folder / "data.hdf5", "r") as f: return f["probabilities"][:] def colorcode_probabilities(probabilities): """Use uint8 to reference classes (we only have 4)""" predicted_classes = np.argmax(probabilities, axis=2).astype(np.uint8) """Do no re-compute max, use argmax result to index""" max_prob = np.take_along_axis(probabilities, predicted_classes[..., None], axis=2).squeeze() """Manual GC, no longer needed""" del probabilities """Use float32 to save memory""" max_prob = max_prob.astype(np.float32) class_colors = np.array([[0, 0, 255], [0, 255, 0], [255, 0, 0], [0, 255, 255]]], dtype=np.uint8) colored_probabilities = class_colors[predicted_classes] """unroll RGB loop""" colored_probabilities[:, :, 0] = colored_probabilities[:, :, 0] * max_prob colored_probabilities[:, :, 1] = colored_probabilities[:, :, 1] * max_prob colored_probabilities[:, :, 2] = colored_probabilities[:, :, 2] * max_prob return colored_probabilities</pre>
	151M	■ 48%		
53%	37M	■ 12%		
33%	38M	■ 12%	133	
	19M	■ 6%		
	14M	■ 4%		
	19M	■ 6%		
	19M	■ 6%		
	19M	■ 6%		

300ms runtime

Line	Time Python	----- <i>native</i>	----- <i>system</i>	
1				<code>import numpy as np</code>
2				<code>import h5py</code>
3				
4				<code>def load_probabilities(folder):</code>
5				<code>with h5py.File(folder / "data.hdf5", "r") as f:</code>
6				<code>return f["probabilities"][:]</code>
7				
8				<code>def colorcode_probabilities(probabilities):</code>
9				<code>"""Use uint8 to reference classes (we only have 4):"""</code>
10	12%	29%		<code>predicted_classes = np.argmax(probabilities, axis=2).astype(np.uint8)</code>
11				<code>"""Do no re-compute max, use argmax result to index:"""</code>
12	2%	9%	4%	<code>max_prob = np.take_along_axis(probabilities, predicted_classes[..., None], axis=2).squeeze()</code>
13				<code>"""Manual GC, no longer needed:"""</code>
14				<code>del probabilities</code>
15				<code>"""Use float32 to save memory:"""</code>
16		1%	16%	<code>max_prob = max_prob.astype(np.float32)</code>
17				<code>class_colors = np.array([[0, 0, 255], [0, 255, 0], [255, 0, 0], [0, 255, 255]]], dtype=np.uint8)</code>
18				<code>colored_probabilities = class_colors[predicted_classes]</code>
19				<code>"""unroll RGB loop:"""</code>
20				<code>colored_probabilities[:, :, 0] = colored_probabilities[:, :, 0] * max_prob</code>
21				<code>colored_probabilities[:, :, 1] = colored_probabilities[:, :, 1] * max_prob</code>
22				<code>colored_probabilities[:, :, 2] = colored_probabilities[:, :, 2] * max_prob</code>
23				<code>return colored_probabilities</code>
--				

Chunking



Chunking

```
def get_chunk_slices(length, chunk_size):
    for start in range(0, length, chunk_size):
        end = min(start + chunk_size, length)
        yield slice(start, end)

def load_probabilities(folder):
    with h5py.File(folder / "data.hdf5", mode="r") as f:
        x_len, y_len = f["probabilities"].shape[:2]
        for x_slice in get_chunk_slices(x_len, CHUNK_SIZE):
            for y_slice in get_chunk_slices(y_len, CHUNK_SIZE):
                yield (x_slice, y_slice), f["probabilities"][x_slice, y_slice]
```

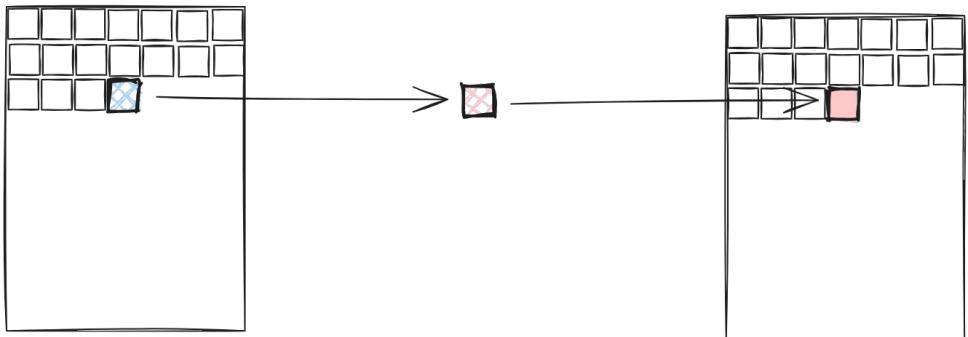
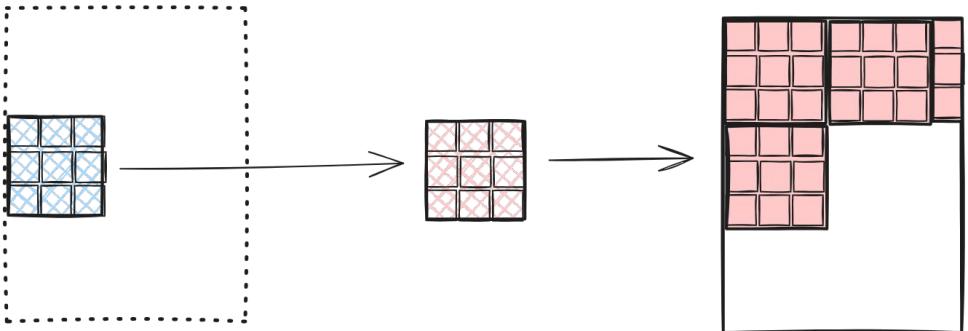
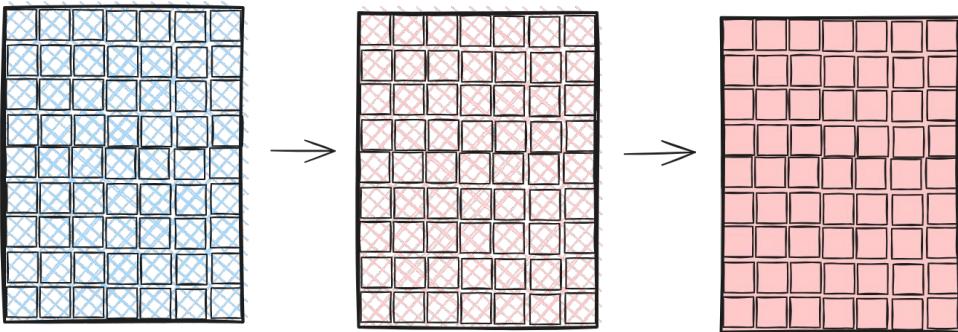
Chunking

```
def load_and_colorcode_probabilities(folder):
    with h5py.File(folder / "data.hdf5", mode="r") as f:
        x_len, y_len = f["probabilities"].shape[:2]
    colored_probabilities = np.zeros(shape=(x_len, y_len, 3), dtype=np.uint8)
    for chunk_slice, chunk in load_probabilities(folder):
        colored_probabilities[chunk_slice] = colorcode_probabilities(chunk)
    return colored_probabilities
```

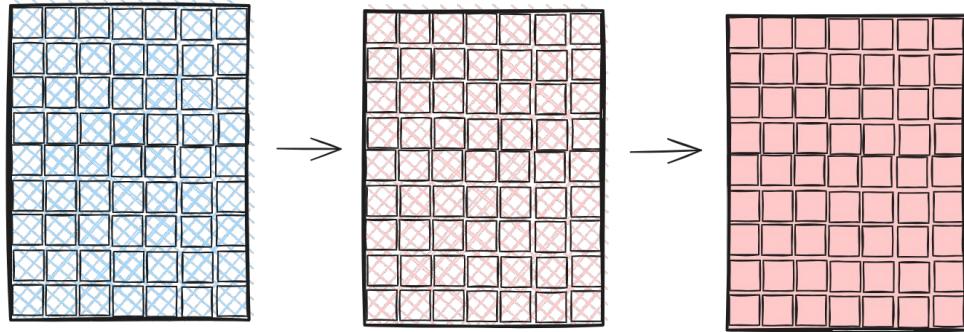
Chunking

	Whole-Array (1920x2560)	Chunked (1024²)
Runtime	300ms	335ms
Memory	211MB	79MB

Element-wise Map with Numba

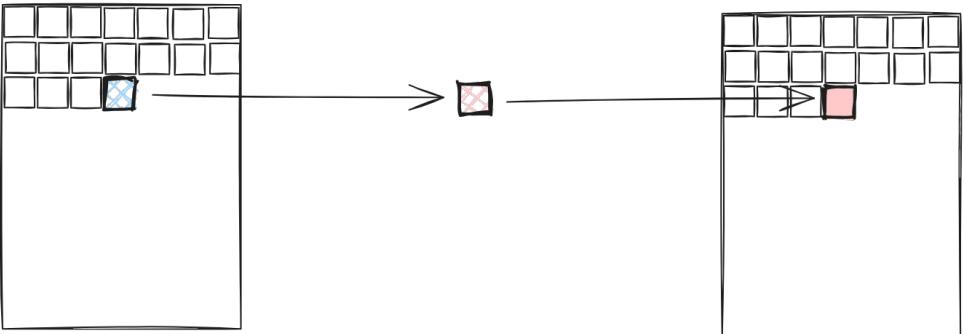
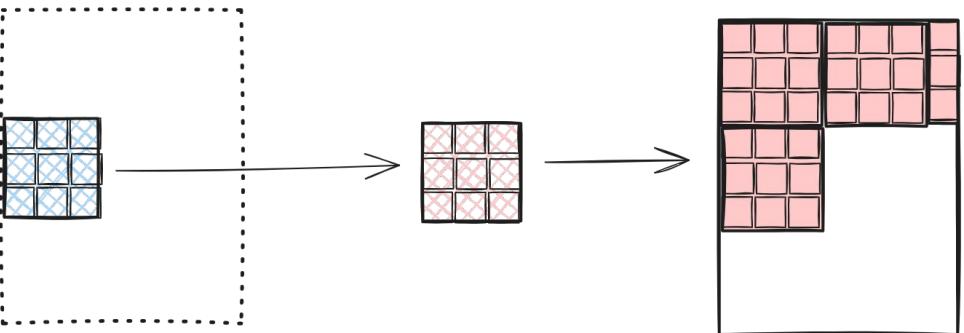


Element-wise Map with Numba



Numba Alternatives

- Pythran
- mypyc
- cython
- ...



Element-wise Map with Numba

```
@numba.njit([numba.uint8[:, :, :](numba.float64[:, :, :, :])])
def colorcode_probabilities(probabilities):
    class_colors = np.array( object: [[0, 0, 255], [0, 255, 0], [255, 0, 0], [0, 255, 255]], dtype=np.uint8)
    x_len, y_len = probabilities.shape[:2]
    colored_probabilities = np.zeros( shape: (x_len, y_len, 3), dtype=np.uint8)
    for x in range(x_len):
        for y in range(y_len):
            class_index = 0
            max_prob = 0
            for i, prob in enumerate(probabilities[x, y]):
                if prob > max_prob:
                    max_prob = prob
                    class_index = i
            colored_probabilities[x, y, 0] = np.uint8(class_colors[class_index, 0] * max_prob)
            colored_probabilities[x, y, 1] = np.uint8(class_colors[class_index, 1] * max_prob)
            colored_probabilities[x, y, 2] = np.uint8(class_colors[class_index, 2] * max_prob)
    return colored_probabilities
```

Element-wise Map with Numba

	numpy	numba
Runtime	300ms	122ms
Memory	211MB	164MB

(non-chunked)

Go More Low-Level

- Custom C++ Extension with `pybind11`
 - Integration of existing c/c++ high-performance code
 - Fine grained control over memory access
 - SIMD
 - Inline Assembly
- Alternatives: PyO3 (rust)

Pybind11 extension

```
py::array_t<uint8_t> colorcode_probabilities(py::array_t<double> probabilities, py::array_t<uint8_t> colors) {
    auto shape = probabilities.shape();
    auto probabilities_ptr = static_cast<double*>(probabilities.request().ptr);
    auto colors_ptr = static_cast<uint8_t*>(colors.request().ptr);
    auto colored_image = py::array_t<uint8_t>({int(shape[0]), int(shape[1]), 3});
    auto colored_image_ptr = static_cast<uint8_t*>(colored_image.request().ptr);
    for (int i = 0; i < shape[0]; i++) {
        for (int j = 0; j < shape[1]; j++) {
            double max_probability = 0.0;
            size_t max_index = 0;
            for (int k = 0; k < shape[2]; k++) {
                auto probability = *probabilities_ptr++;
                if (probability > max_probability) {
                    max_probability = probability;
                    max_index = k;
                }
            }
            auto class_color_ptr = colors_ptr + (max_index * 3);
            *colored_image_ptr++ = static_cast<uint8_t>(*class_color_ptr++ * max_probability);
            *colored_image_ptr++ = static_cast<uint8_t>(*class_color_ptr++ * max_probability);
            *colored_image_ptr++ = static_cast<uint8_t>(*class_color_ptr++ * max_probability);
        }
    }
    return colored_image;
}
```

Pybind11 extension

	numpy	numba	pybind11
Runtime	300ms	122ms	119ms
Memory	211MB	164MB	164MB

(non-chunked)

Repeat

- Profile
- Efficient IO
- Vectorization
- Reduce Memory Usage
- Go More Low-Level

Optimize



15s – 800MB



150ms – 80MB

Now we can scale, e.g. with

- threading
- multiprocessing
- joblib, Spark, Dask, Ray, ...



Further reading:

- codical.org
- pythonspeed.com/datascience
- calmcode.io



Thanks!

code & slides

github.com/jstriebel/profile-optimize-repeat



Valentin Nieper | valentin.nieper@scalableminds.com

Jonathan Striebel | jonathan@aignostics.com  @jostriebel