

Åbo Akademi

From Python Scripting to Parallel Spatial Modeling

Jesús Carabaño Bravo

jcaraban@abo.fi

PhD Candidate at **ÅBO AKADEMI**

Faculty of Science and Engineering

Introduction

We are building a **prototype** of a spatial **framework** / library / module that uses **compilers** techniques to automatically **optimize** sequential raster scripts in e.g. python and executes scripts in **parallel** on e.g. GPUs so that we can handle **very large** datasets

*Looks like raster processing in **ArcPy** / **Matlab** / **Numpy***

Applications: land use, hydrology, air quality, land erosion, predictive analysis, geomorphology, ecology

Script 1: urban development

```
1 from map import * ## "Parallel Map Algebra" package
2
3 a = 6.4640 # Constant coefficient
4 b1 = 43.5404 # Elevation coefficient
5 b2 = 1.9150 # Slope coefficient
6 b3 = 41.3441 # Distance to city centers coefficients
7 b4 = 12.5878 # Distance to transportations coefficient
8 b5 = [0,0,-9.865,-8.746,-9.268,-8.032,-9.169,-8.942,-9.45]
9 # {water,urban,barren,forest,shrub,woody,herb,crop,wetlad}
10 d = 5 # dispersion parameter
11 q = 16000 # max cells to become urban per year
12
13 x1 = read('dem') # elevation layer
14 x2 = read('slope') # slope layer
15 x3 = read('center') # distance to centers layer
16 x4 = read('transp') # distance to transportations layer
17 x5 = read('landuse') # land use layer
18 e = read('excl') # exclusion layer (e.g. water bodies)
19 s = read('urban') # initial state: urban / not-urban
20 N = 50 # years of simulation i.e. time steps
21
22 for i in range(N) :
23     z = a + b1*x1 + b2*x2 + b3*x3 + b4*x4 + pick(x5,b5)
24     pg = exp(z) / (1 + exp(z))
25     pc = pg * !e * !s * focalSum(s) / (3*3-1)
26     pd = pc * exp(-d * (1 - pc / zonalMax(pc)))
27     ps = q * pd / zonalSum(pd)
28     s = s || ps > rand()
29
30 write(s, 'output')
```

Ref: Wu 2002 "Calibration of stochastic cellular automata: the application to rural-urban land conversions"

Script 1: urban development

```
1 import urban # imports 'urban()' function, containing Listing 1
2
3 prob = zeros() # urban probability map
4 M = 1000      # Monte Carlo iterations
5
6 for i in range(0,M) :      # Monte Carlo method
7     prob = prob + urban() # urban() returns the urban layer
8     prob = prob / M       # urban() ∈ {0,1} ==> prob ∈ [0,1]
9
10 write(prob, 'output')
```

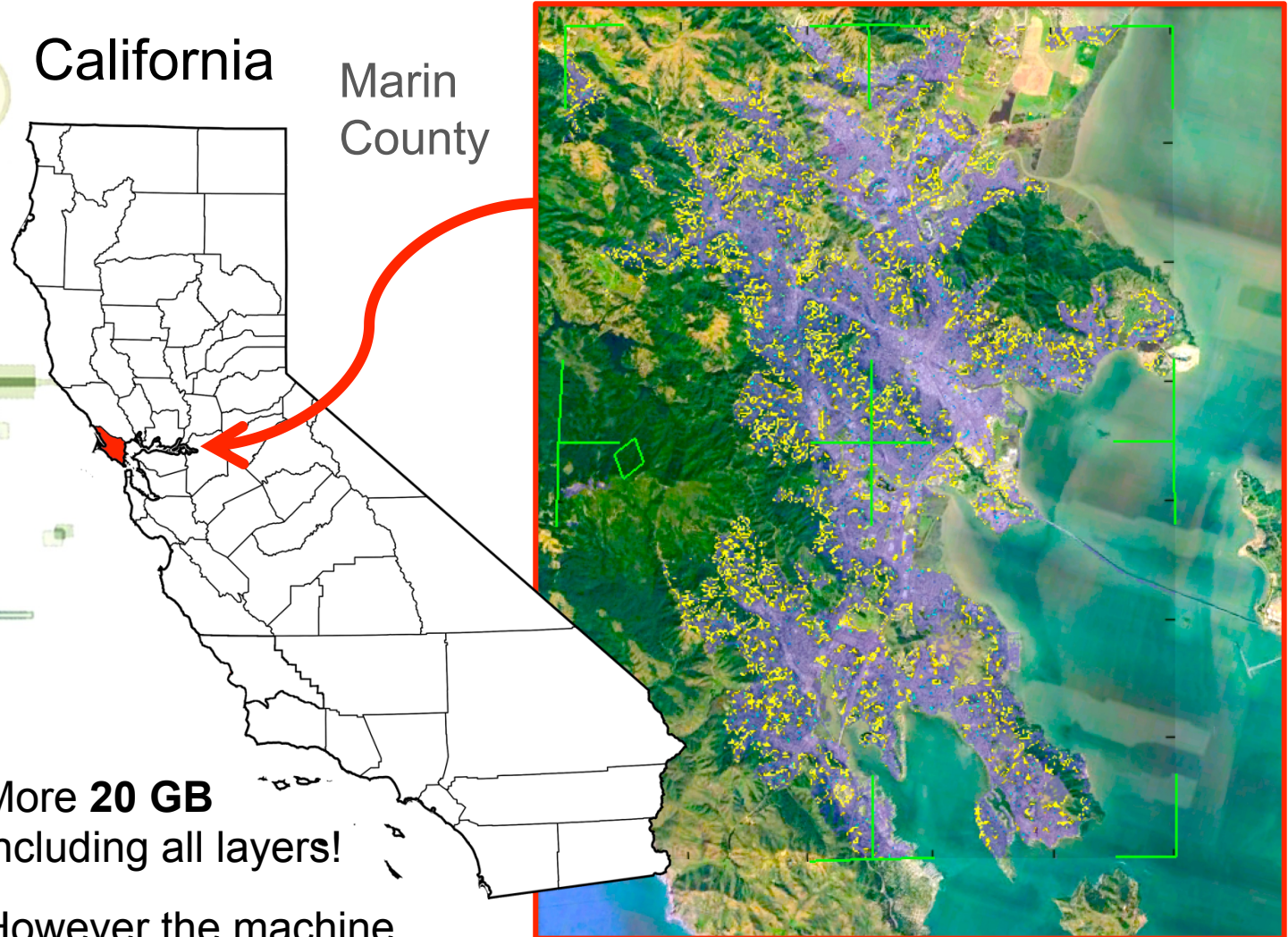
Optimizations	Baseline	GPU	GPU + Loc	GPU + Loc + Act
Execution	20 min	10 min	3 min	1 min 10 s
Speed Up	1x	2x	6.66x	17.14x
<i>Baseline = CPU SIMD Threaded Loc = Locality Opt. Act = Only Active Blocks</i>				
Monte Carlo	1 iteration	10 iter.	100 iter.	1000 iter.
Execution	1 min 14 s	12 min	121 min	1214 min

Ref 1 = 32s with 64 GPUs, We = 70s with just 1 GPU

Ref 1: Guan 2016 “A hybrid parallel cellular automata model for urban growth simulation over GPU/CPU...”

Ref 2: Guan 2014 “pRPL 2.0: Improving the Parallel Raster Processing Library,” Trans. GIS, vol. 18, ...

Script 1: urban development



California

Marin
County

More **20 GB**
including all layers!

However the machine
only has **16 GB** of RAM

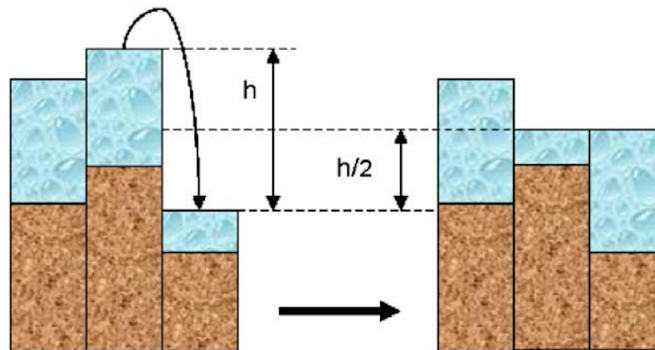
Script 2: flooding model

```
1 from map import * ## "Parallel Map Algebra" pac
2
3 h = read('dem')      # digital elevation layer
4 w = read('water')    # water depth layer
5 i = read('inflow')  # inlets inflow layer
6 o = read('outflow') # outlets outflow layer
7 N = 1000             # number of time steps
8
9 def swap(x,i,j) :
10     x[i], x[j] = min(x[i],x[j]), max(x[i],x[j])
11
12 def netsort5(x) :
13     swap(x,0,1); swap(x,2,3); swap(x,0,2)
14     swap(x,3,4); swap(x,0,3); swap(x,1,3)
15     swap(x,2,4); swap(x,1,4); swap(x,1,2)
16
17 def avglevel(w,h,x) :
18     netsort5(x) # ascending order
19     s = w+x[0] # sum variable
20     n = 1      # count variable
21     for i in range(1,5) :
22         b = (s >= x[i]*i)
23         s += b*x[i]
24         n += b
25     return s / n
26
27 def gather(w,h) :
28     x = [0]*5 # neighborhood (NBH)
29     x[0] = h # central cell
30     x[1] = h[0,-1] + w[0,-1]
31     x[2] = h[-1,0] + w[-1,0]
32     x[3] = h[+1,0] + w[+1,0]
33     x[4] = h[0,+1] + w[0,+1]
34     return avglevel(w,h,x)
35
36 def distri(w,h,l) :
37     wh = w+h # prev water level
38     c = max(0, l[0,-1] - wh)
39     c += max(0, l[-1,0] - wh)
40     c += max(0, l[+1,0] - wh)
41     c += max(0, l[0,+1] - wh)
42     c += max(h, l) - wh
43     cwh = max(c + wh, h)
44     return cwh - h
45
46 for j in range(0,N) :
47     w = w + i # fill inlet
48     l = gather(w,h) # gather avg
49     w = distri(w,h,l) # distribute
50     w = max(w-o,0) # drain water
51
52 write(w,'output')
```

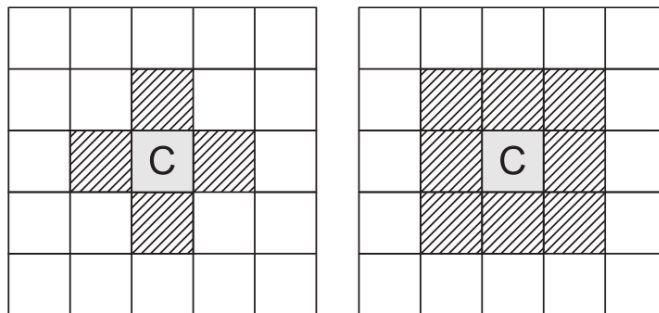
Ref: S. Di Gregorio and R. Serra, "An empirical method for modelling and simulating some complex macroscopic phenomena by cellular automata"

Ref: P. Topa, "Cellular Automata Model Tuned for Efficient Computation on GPU with Global Memory Cache"

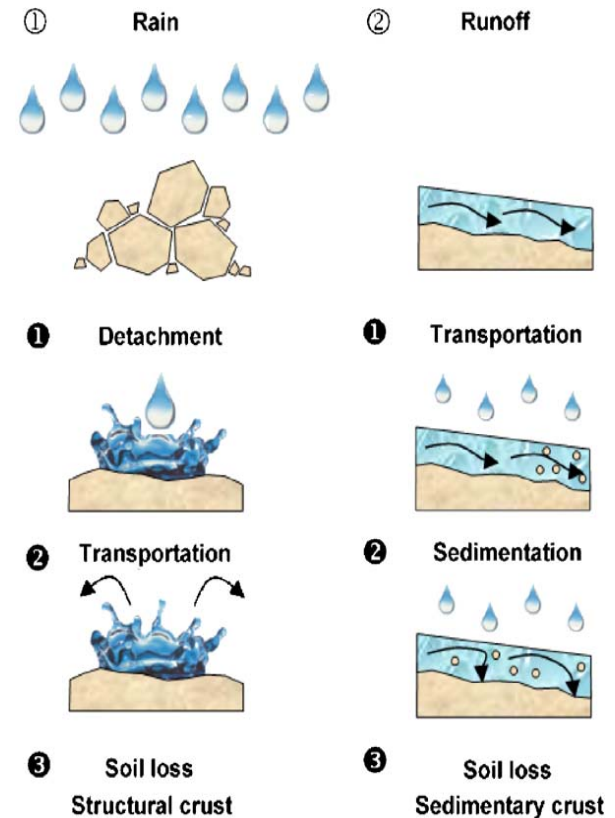
Script 2: flooding model



1-dimensional water balance



2-dimensional neighborhood



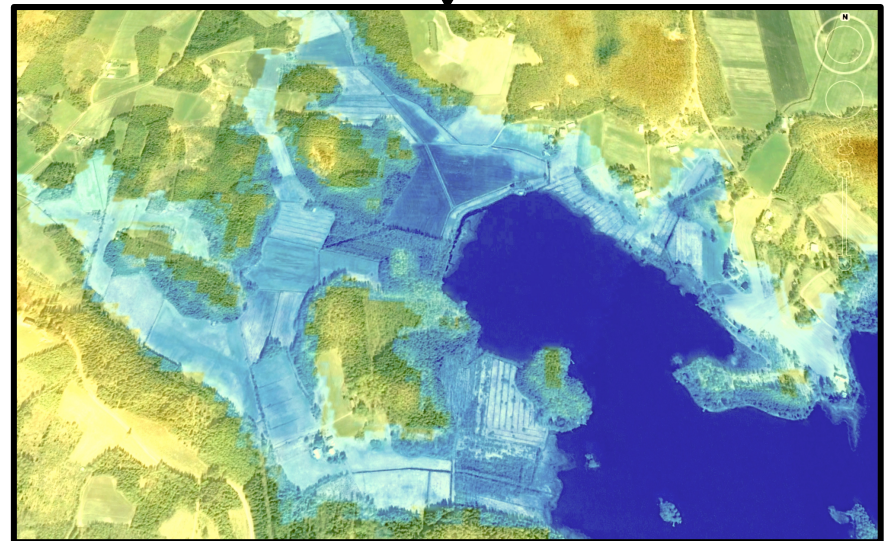
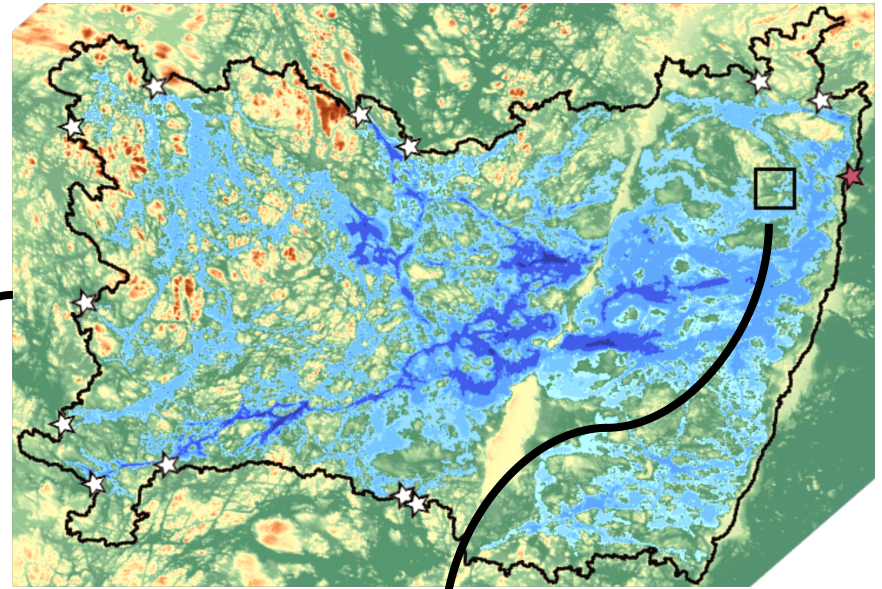
More advanced erosion model

Ref: Valette 2016 "SoDA project: A simulation of soil surface degradation by rainfall"

Ref: D. D'Ambrosio 2001 "A Cellular Automata model for soil erosion by water"

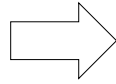
Script 2: flooding model

Saimaa Basin
Finland

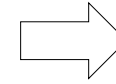


More than magic

Python
Script



Compiler
Magic



Parallel
Speedup

```
27 def gather(w,h) :
28     x = [0]*5 # neighborhood (NDB)
29     x[0] = b # central cell
30     x[1] = b[0,-1] + w[0,-1]
31     x[2] = b[-1,0] + w[-1,0]
32     x[3] = b[+1,0] + w[+1,0]
33     x[4] = b[0,+1] + w[0,+1]
34     return avglevel(w,h,x)
35
36 def distri(w,h,l) :
37     wh = w*h # prev water level
38     c = max(0, l[0,-1] - wh)
39     c += max(0, l[-1,0] - wh)
40     c += max(0, l[+1,0] - wh)
41     c += max(0, l[0,+1] - wh)
42     c += max(h, l) - wh
43     cwh = max(c + wh, h)
44     return cwh - h
45
46 for j in range(0,N) :
47     w = w + i # fill inlet
48     l = gather(w,h) # gather avg
49     w = distri(w,h,l) # distribute
50     w = max(w-o,0) # drain water
51
52 write(w, 'output')
```



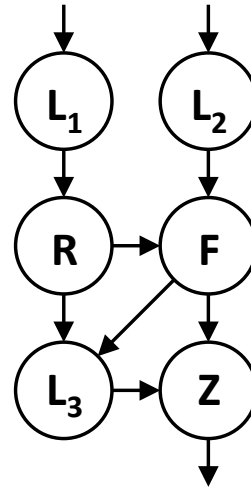
Compiler Techniques

```

IN1 = read(input1)
IN2 = read(input2)
L1 = LocalOp(IN1)
L2 = LocalOp(IN2)
R = RadialOp(L1)
F = FocalOp(R,L2)
L3 = LocalOp(R,F)
Z = ZonalOp(L3,F)
write(Z,output)
    
```

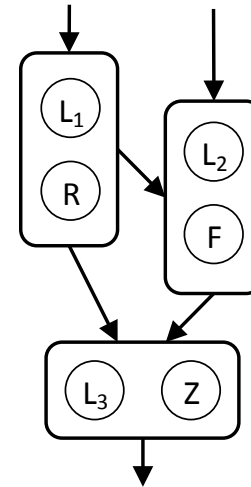
a) Python Script

Symbolic



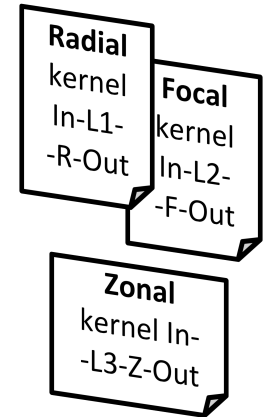
b) Dependency Graph

Fusion

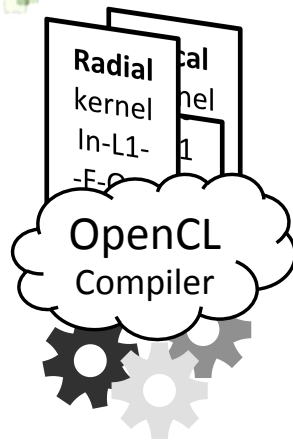


c) Grouped Graph

Generation

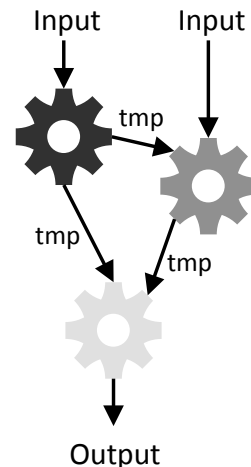


d) GPU Code



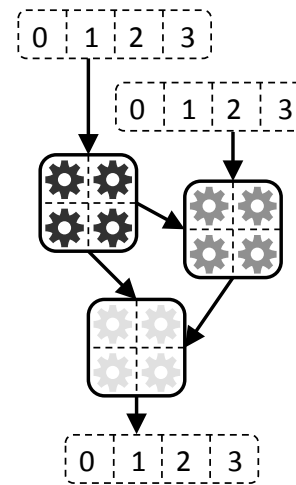
e) Compilation

Executable



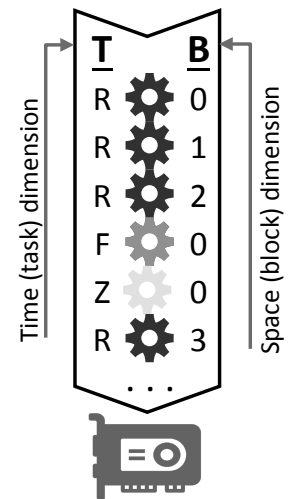
f) Tasks & Dataflow

Decomposition



g) Blocks, Jobs

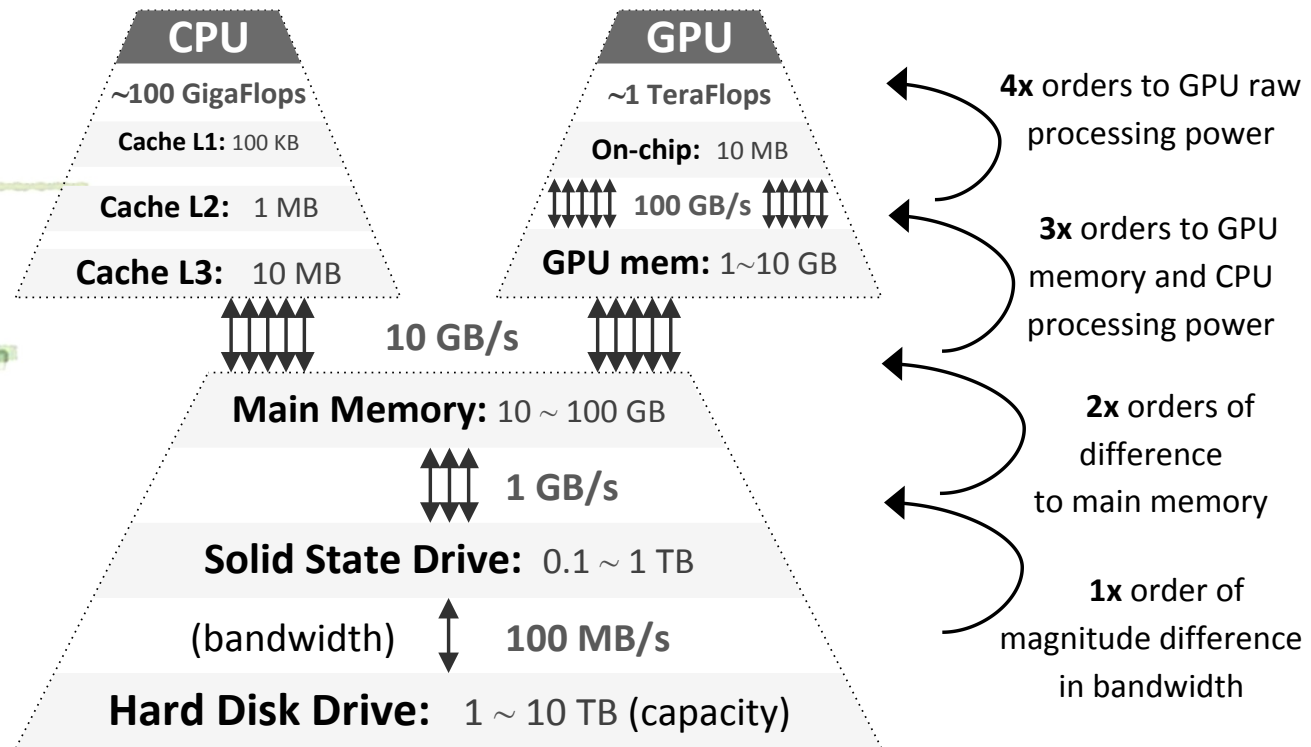
Reordering



h) Scheduling

Struggling for locality

*The #1 optimization is **locality**. Avoids idle CPU cycles waiting for the data*



GRASS, QGIS, ArcGIS don't exploit locality!

Summary

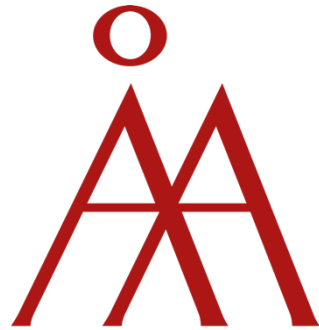
A framework that automatically optimizes raster python scripts and executes them in parallel.

*Looks like **ArcPy** / **Matlab** / **Numpy***

Applications: land use, hydrology, air quality, land erosion, predictive analysis, geomorphology, ecology

Consequences

- You might not need supercomputing power
- because your machine is still underutilized
- and workstations are easier to work with!



Åbo Akademi

Thanks for your time!

Jesús Carabaño Bravo

jcaraban@abo.fi

PhD Candidate at **ÅBO AKADEMI**

Faculty of Science and Engineering