



# Managing Geospatial Big Data: Lessons Learned and Future Perspectives

Marcos Vaz Salles

Associate Professor, University of Copenhagen (DIKU)

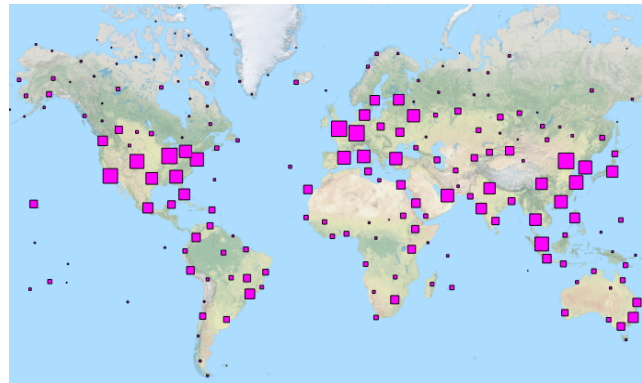
# Geospatial Applications

- Notion of geographical space fundamental to human activities and the environment
- Applications with spatial characteristics among the most exciting in computing

<https://media.licdn.com/mpr/mpr/jc/AAEAAQAAAAAAAAAX8AAAAJDEyMjY2Zjc3LWIyOWMtNDY2YS1hZGE2LWRjNDU3MDg2OTE0ZQ.jpg>



Location-based services,  
augmented reality



Data visualizations,  
e.g., maps

[https://en.wikipedia.org/wiki/Spatial\\_analysis#/media/File:Snow-cholera-map.jpg](https://en.wikipedia.org/wiki/Spatial_analysis#/media/File:Snow-cholera-map.jpg)

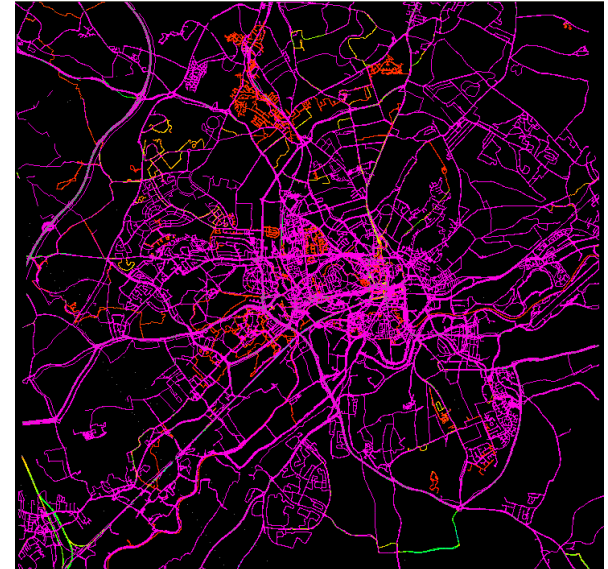


Spatial analysis

# Challenges in Geospatial Applications

- **Big variation in data formats and volume**
  - Some data is “cheap” to obtain, other is “expensive”
  - Examples: Drone images, GPS traces, satellite data vs. visual ranking of breeds by humans

[https://en.wikipedia.org/wiki/Unmanned\\_aerial\\_vehicle#/media/File:Intersect\\_UAV\\_B\\_3.1.png](https://en.wikipedia.org/wiki/Unmanned_aerial_vehicle#/media/File:Intersect_UAV_B_3.1.png)



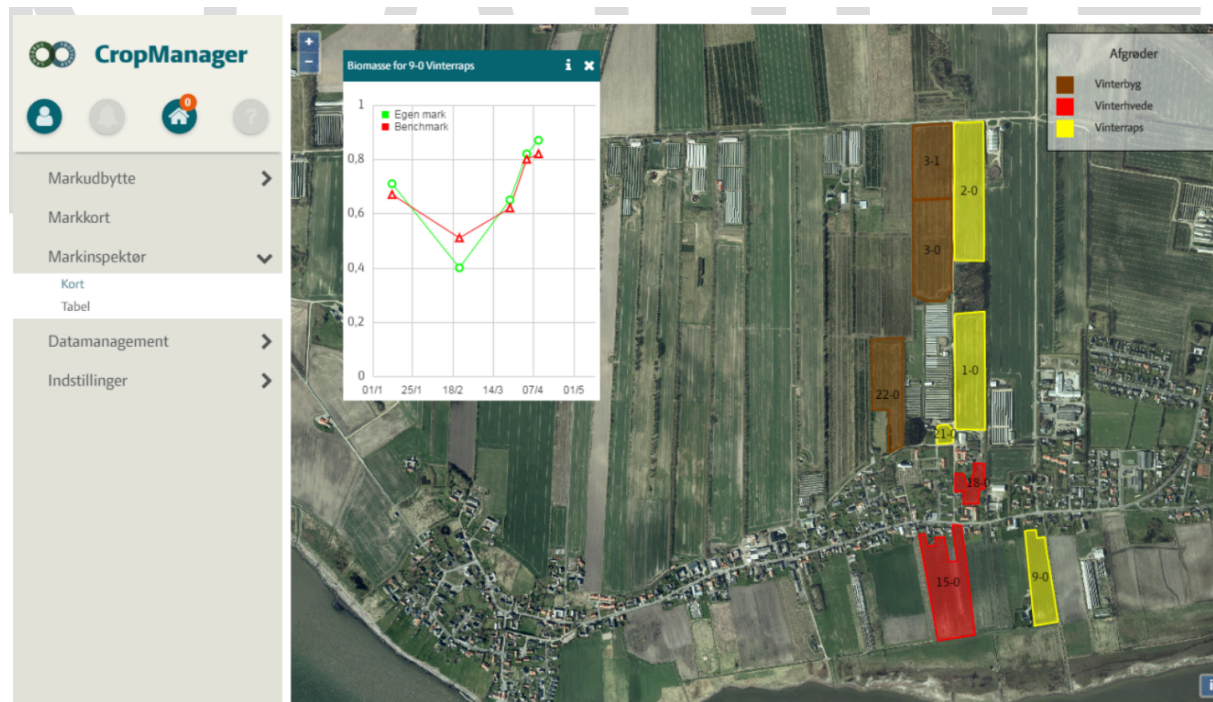
[http://wiki.openstreetmap.org/w/images/e/e7/Nottingham\\_gps\\_traces\\_ex\\_osm\\_20110105.png](http://wiki.openstreetmap.org/w/images/e/e7/Nottingham_gps_traces_ex_osm_20110105.png)



<http://futurecropping.eng.au.dk/maps/204>

# Challenges in Geospatial Applications

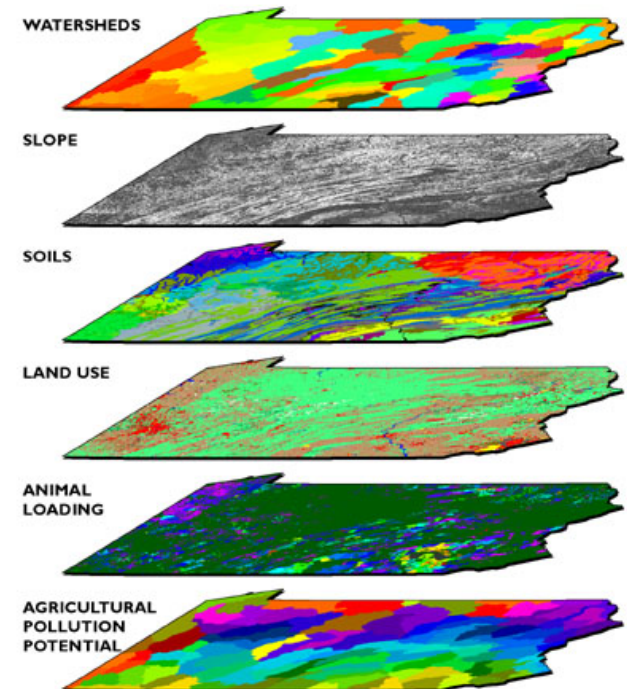
- **Large amount of users and potentially complex simultaneous requests**
  - Popular datasets need to be serviced to many users, and transformed by different programs
  - Examples: Google Maps & shortest paths, Future Cropping data platform & analysis services



NDVI benchmarking in SEGES Crop Manager, tech transfer from Future Cropping project

# Challenges in Geospatial Applications

- **Much labor needed to derive knowledge from varied data**
  - “Expensive” data can be too small or too noisy for phenomenon studied, not obvious how to leverage “cheap” data
  - Examples: Soil samples vs. satellite or land use data



[https://www.e-education.psu.edu/natureofgeoinfo/c9\\_p6.html](https://www.e-education.psu.edu/natureofgeoinfo/c9_p6.html)

# Lessons From Managing Geospatial Data

- **Challenge: Big variation in data formats and volume**
  - Lesson 1: “Cheap” vs. “expensive” data
  - Lesson 2: The rise of standardization, open-source software, and large geospatial datasets
- **Challenge: Large amount of users and potentially complex simultaneous requests**
  - Lesson 3: From software to services
  - Lesson 4: Telemetry turns behavior into data
- **Challenge: Much labor needed to derive knowledge from varied data**
  - Lesson 5: Embed intelligence in services



# Lesson 1: “Cheap” vs. “Expensive” Data



# What makes data “expensive”?

<http://www.agronomy.k-state.edu/services/soiltesting/farmer-services/soil-analysis/index.html>

- **Lack of automated sensing technology**
  - For example, soil samples or phenotype annotations require human labor *for each sample*
- **Lack of data description (metadata) or data quality controls**
  - For example, drone images can turn out to be very *noisy data* due to measurement errors, including alignment, color filters, variability in cloud cover, variability in RGB profiles across drones



[https://www.nordgen.org/ngdoc/plants/ppp\\_sekr/PPP\\_Basic\\_Documents/Basic\\_documents/ppp\\_promoting\\_nordic\\_plant\\_breeding\\_for\\_the\\_future.pdf](https://www.nordgen.org/ngdoc/plants/ppp_sekr/PPP_Basic_Documents/Basic_documents/ppp_promoting_nordic_plant_breeding_for_the_future.pdf)



# What makes data “expensive”?

- **Lack of data-centric organizational culture, tools, and technology**
  - For example, top management does not see data as first-class entity in business, or there is no data team or platform established
- **Lack of possibilities to externalize data management costs**
  - For example, there may be no service providers in the given area, or data may be strategic differentiator

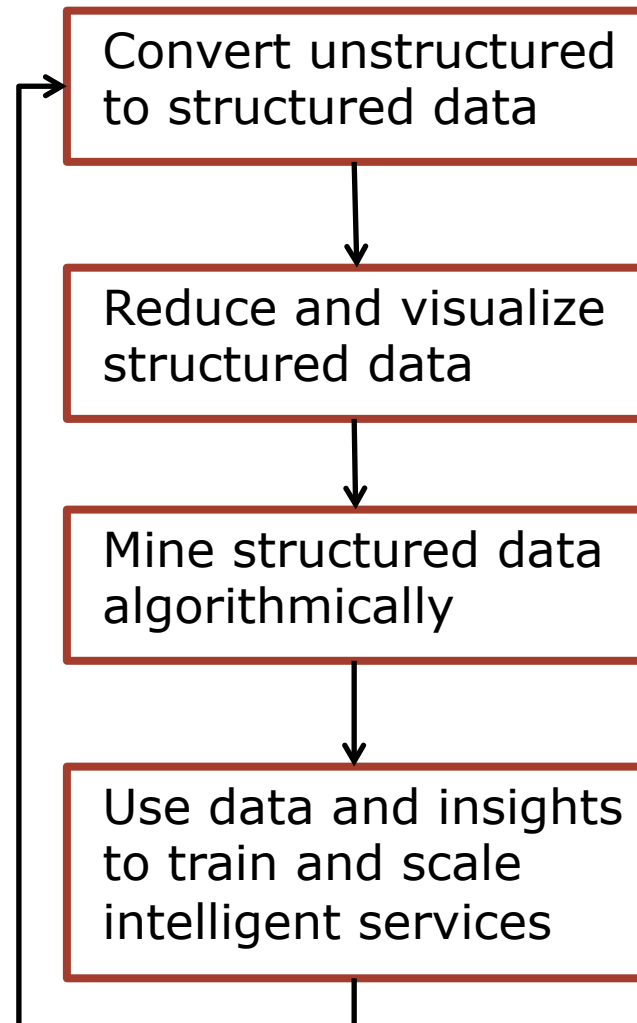
<https://dribbble.com/shots/2055315-We-Love-Data>



<http://gst.dk/>

Sentinel-2 satellite, image by Rama,  
[https://upload.wikimedia.org/wikipedia/commons/3/3d/Sentinel\\_2-IMG\\_5873-white\\_%28crop%29.jpg](https://upload.wikimedia.org/wikipedia/commons/3/3d/Sentinel_2-IMG_5873-white_%28crop%29.jpg)

## Becoming Data-Centric



- **Every step is important:** Partial achievement is possible
- Achieving the whole loop requires teams of both **data engineers** and **data scientists**

# The challenges of exploiting “cheap” data

- Data as **intangible asset**
  - Investment needed to deliver value!
- Getting to **structure**: Goal is to represent data as **tables** (preferred) or **matrices**
  - **Non-trivial transformations** to structure data, e.g., how to think of free text or images as tables or matrices?
  - **Heterogeneity** in representation of data across different databases in different table formats (**schema-level**) or of same data in different sources with different attributes (**instance-level**)
  - **Errors** in data leading to the need for data quality procedures and **data cleaning**



# Lesson 2: The Rise of Standardization, Open-Source Software, and Large Geospatial Datasets



# Standardization in Geospatial Data

- **Open Geospatial Consortium (OGC)**

- Standards body for geospatial data
- Work on data formats, e.g., well known text (WKT), NetCDF
- Work on protocols, e.g., WFS, WCS, WMS, WPS

<http://www.opengeospatial.org/>



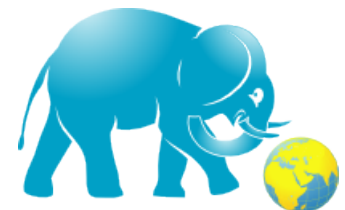
- **Open-source software**

- *Spatial application server and CMS*: GeoServer, GeoNode
- *Spatial relational database*: PostgreSQL/PostGIS
- *Raster analytics*: Rasdaman
- *Spatial Big Data*: Spatial Hadoop, Simba, GeoMesa, GeoWave

<http://www.osgeo.org/>



<http://spatialhadoop.cs.umn.edu/>



<http://www.rasdaman.org/>

**rasdaman**  
raster data manager



## Standardization helps in structuring data

- **Common data model** eases non-trivial transformations, provides way to leverage previous efforts
- **Open-source software** reduce data heterogeneity at both schema and instance levels
- **Standard formats** allow for large, high-quality datasets to be curated and shared
- **BUT...**
  - Standardization tends to work when data targeted is supposed to become *commonly used across industry*
  - *Spreads costs* of integrating data among participants
  - Hard to achieve when data provides *proprietary competitive advantage* to organizations



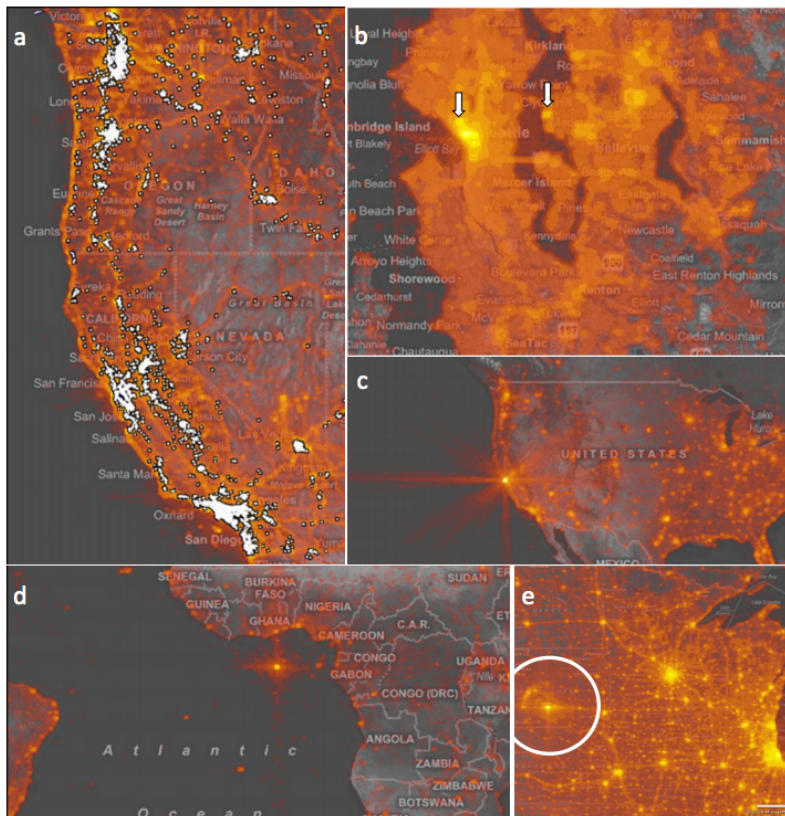
- Standards and open-source software pave the way to large datasets
- **Concern: What if the data does not fit my spreadsheet?**

Data provided by Rasmus L. Hjortshøj at Sejet



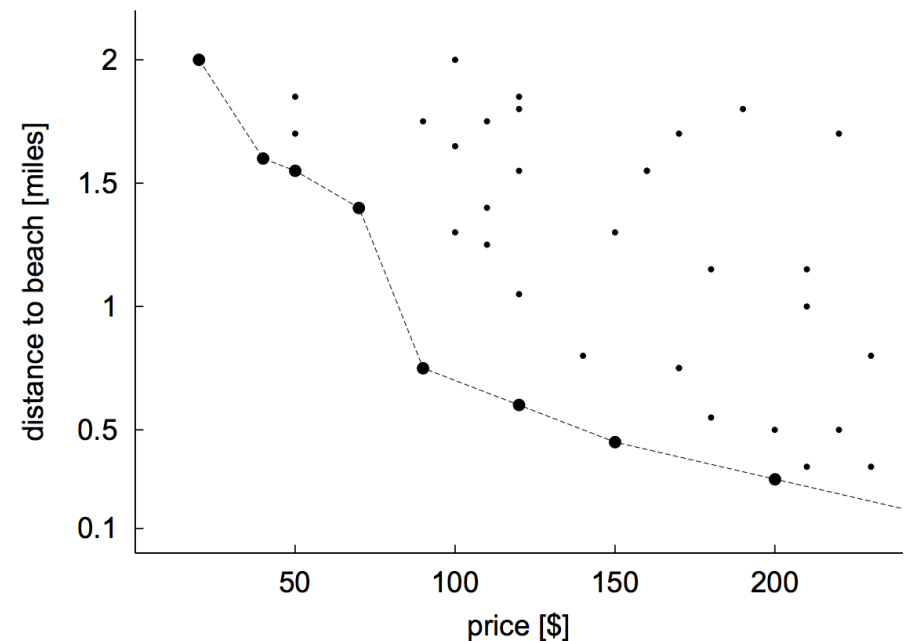
# Data Reduction Methods

- **Data aggregation,**  
e.g., heatmaps



"Hotmap: Looking at Geographic Attention",  
Danyel Fisher, Microsoft Research (2007)

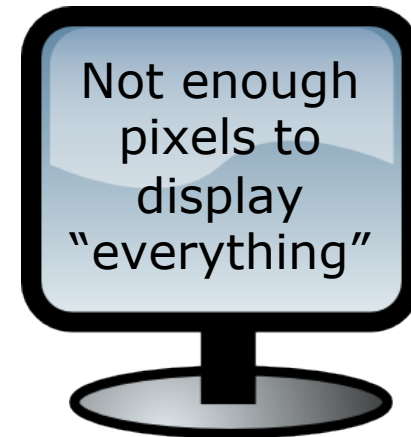
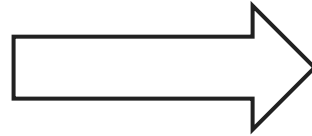
- **Data selection,**  
e.g., skyline,  
cartographic  
generalization



"The Skyline Operator", Stephan Börzsönyi,  
Donald Kossmann, Konrad Stocker, ICDE 2001

## Research Highlight: Declarative Cartography

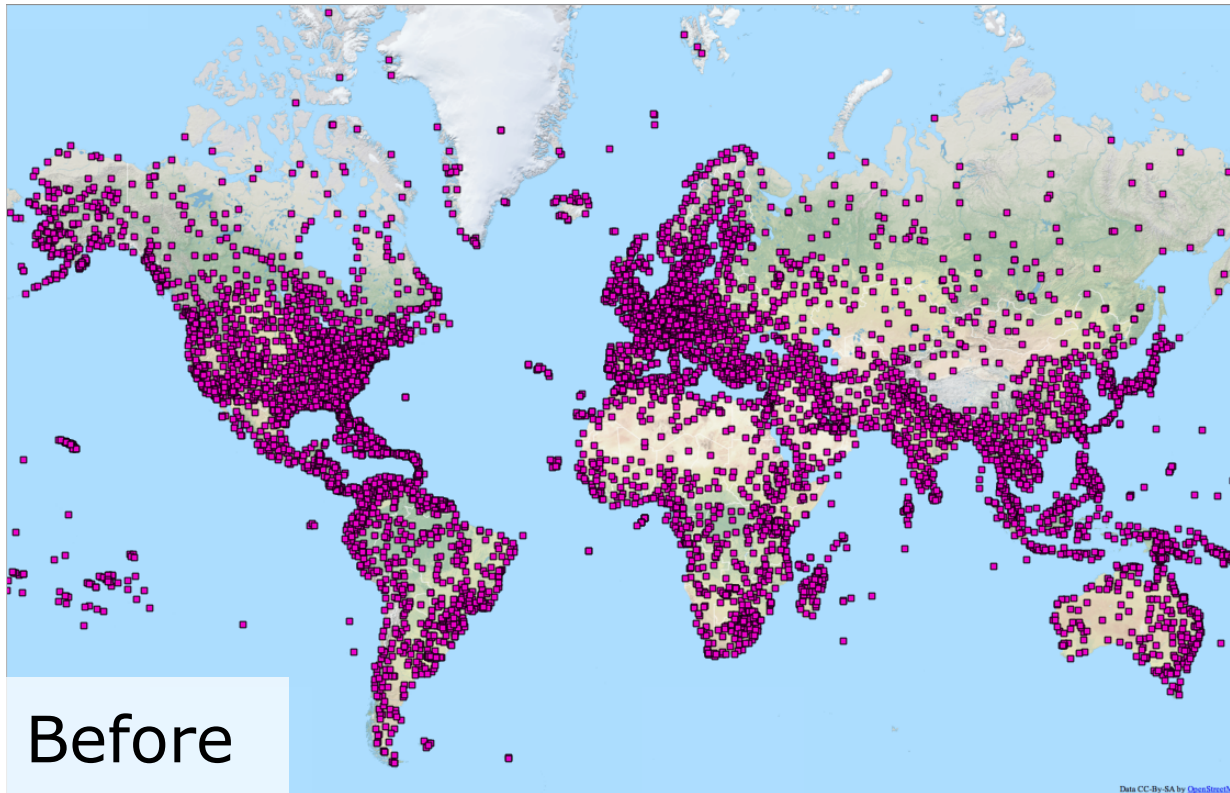
Adapting data to scale of visualization medium



Work done in collaboration with P. K. Kefaloukos  
and M. Zachariasen, results in ICDE 2014

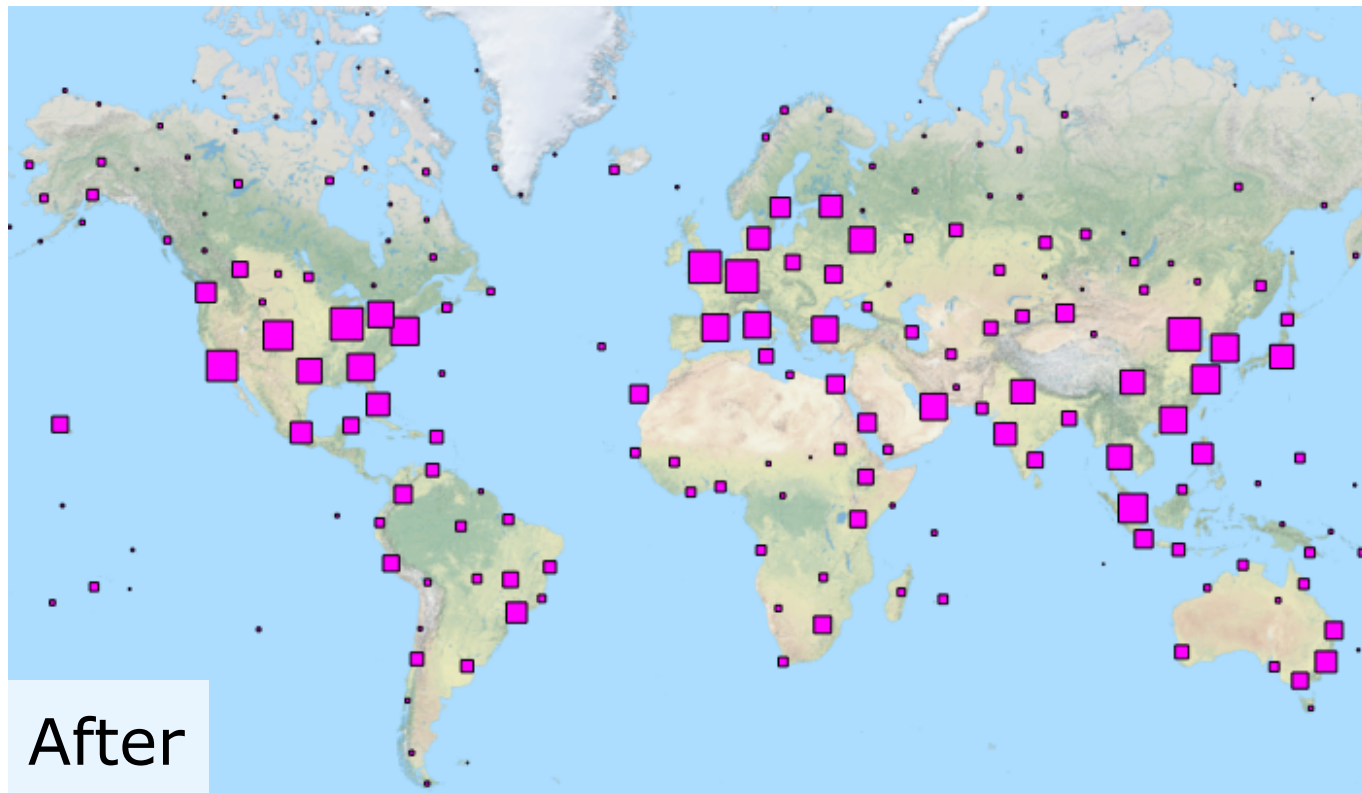


## Example: Selecting Airports



- Too much information → illegible map
- Not clear how to deal with zooming
- Not obvious how to pick objects to display

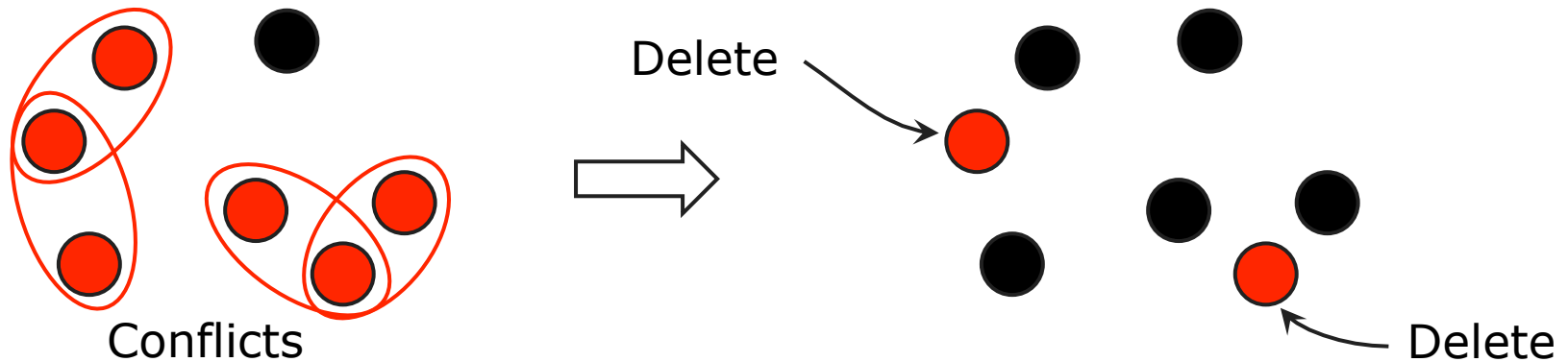
## Example: Selecting Airports



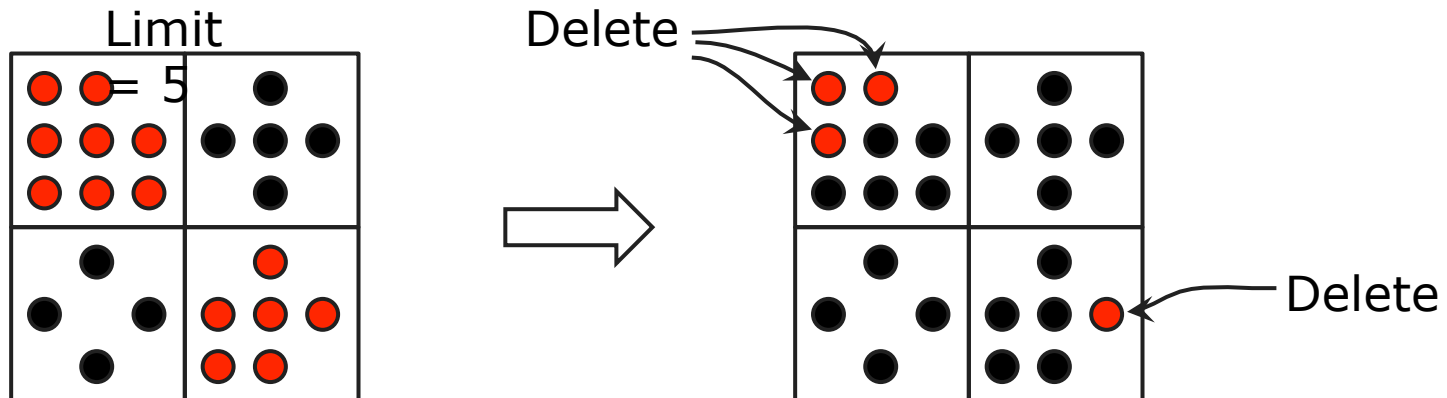
How can we get from "before" to "after"?

# Cartographic Constraints

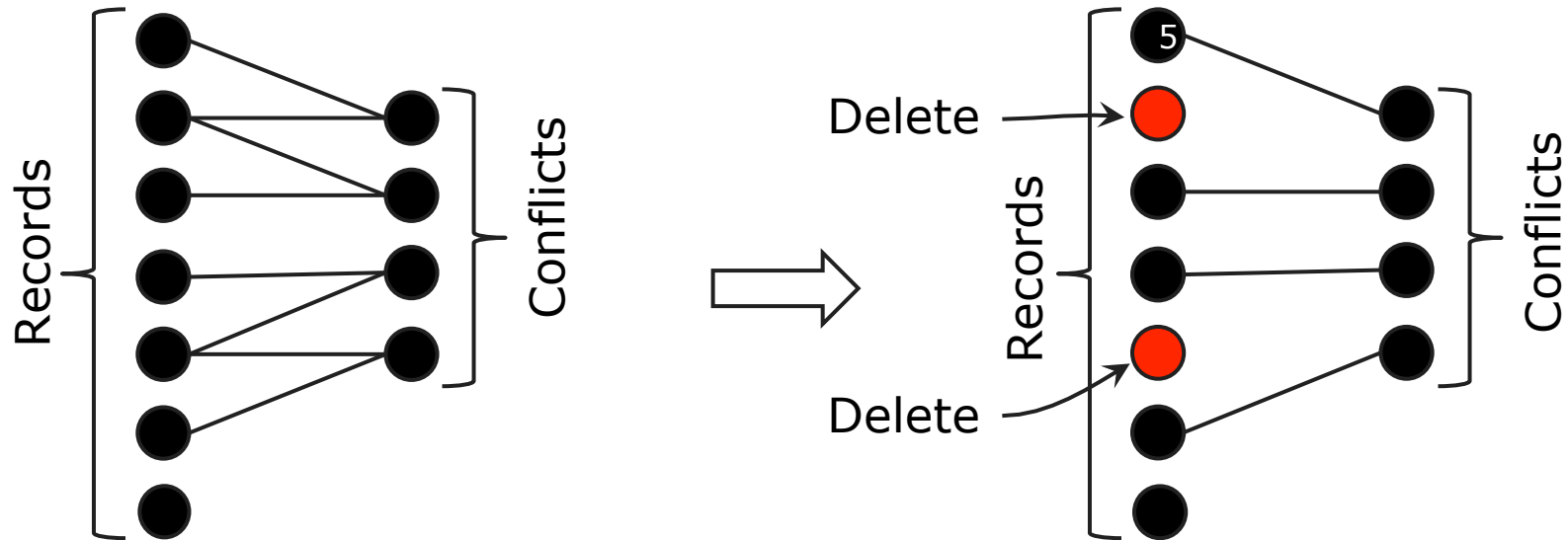
**Proximity constraint:** minimum distance between records (measured in pixels on screen)



**Visibility:** maximum records per unit area (within a map "tile")



# Optimization Problem



- Cartographic constraints and record importance lead to optimization problem
- Delete minimum weight cover
- Set multicover problem  $\rightarrow$  NP-Hard

# Declarative Cartography



## GENERALIZE

airports TO airports\_zoom

WITH ID airport\_id

WITH GEOMETRY wkb\_geometry

AT 18 ZOOM LEVELS

## WEIGH BY

num\_departures

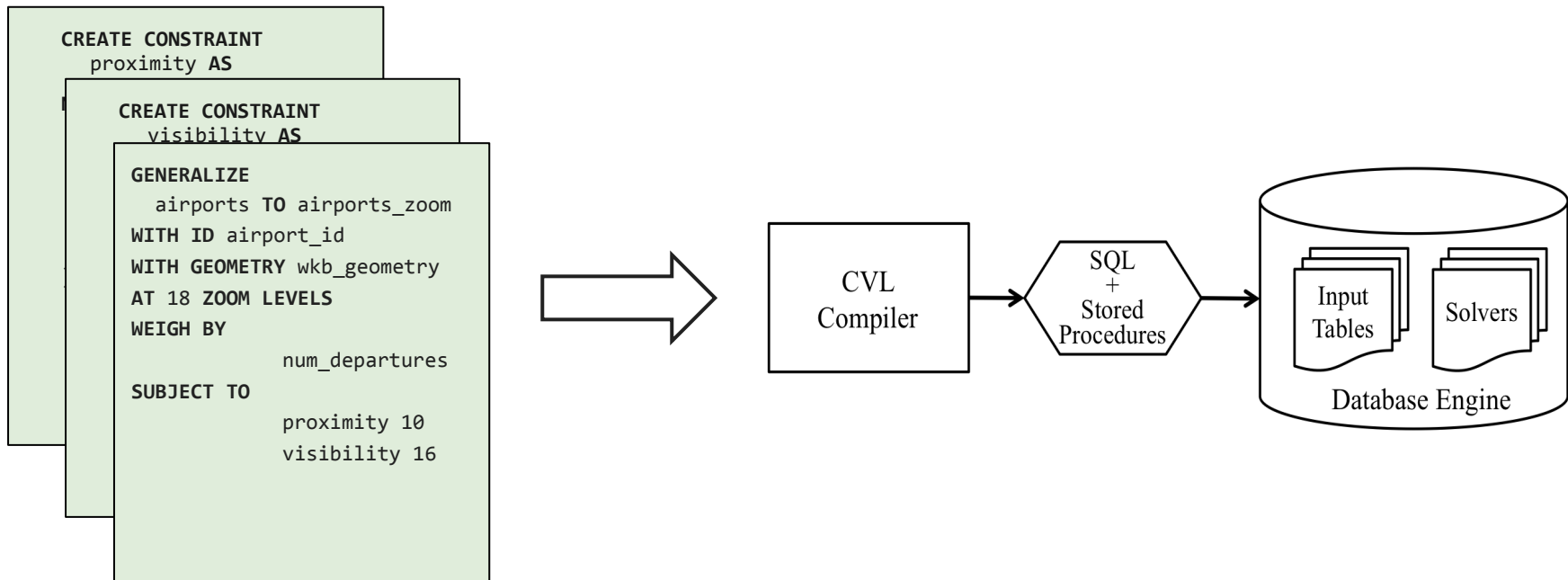
## SUBJECT TO

proximity 10

visibility 16

- Creating maps is the job of data-journalists, bloggers, high-level programmers... not mathematicians
- **Cartographic Visualization Language (CVL):** transforms input data into zoomable data
- Constraints expressed in SQL

# Compiling CVL into SQL



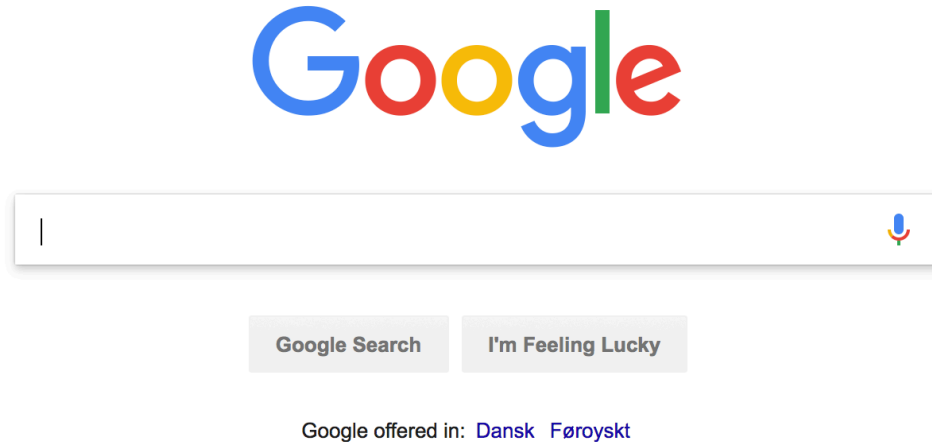
- **Leverage database theory and technology to compute generalization inside DB (in-situ)**
- Our prototype: PostgreSQL + PostGIS + Python + CVXOPT
  - Could be any database, e.g., a parallel one!

# Lesson 3: From Software to Services



# Did you ever install Google?

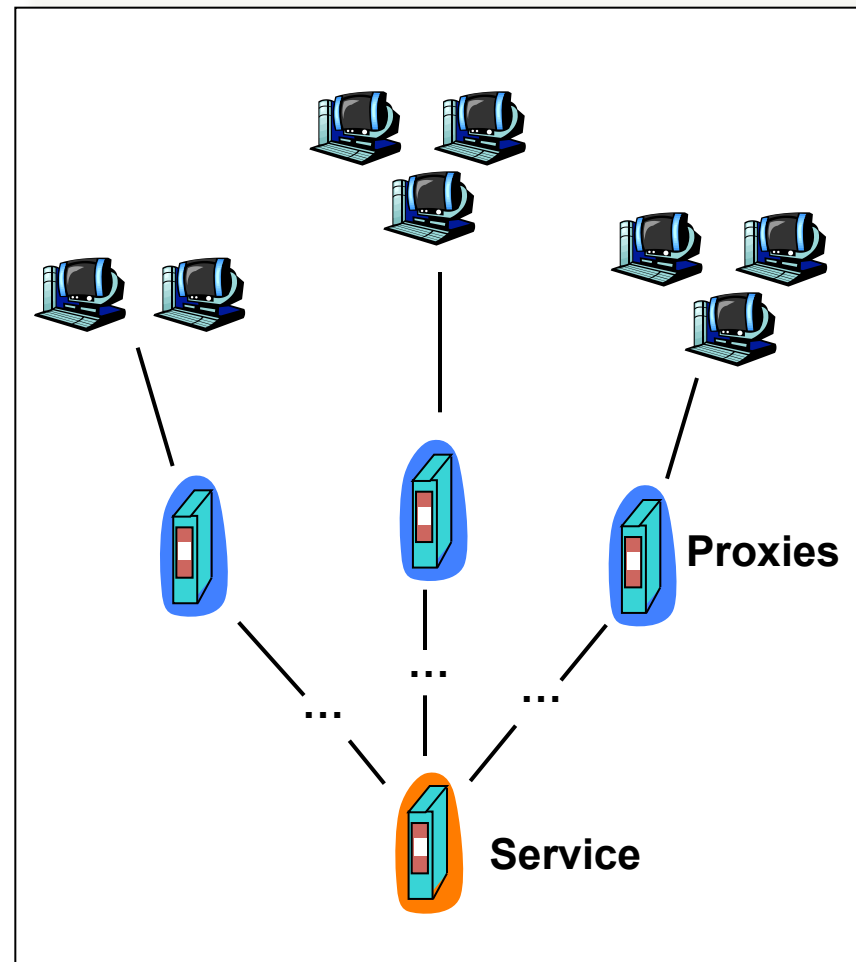
[https://www.google.dk/?gws\\_rd=cr&dcr=0&ei=AxsIWu-8MouE6AS-3YKgDA](https://www.google.dk/?gws_rd=cr&dcr=0&ei=AxsIWu-8MouE6AS-3YKgDA)



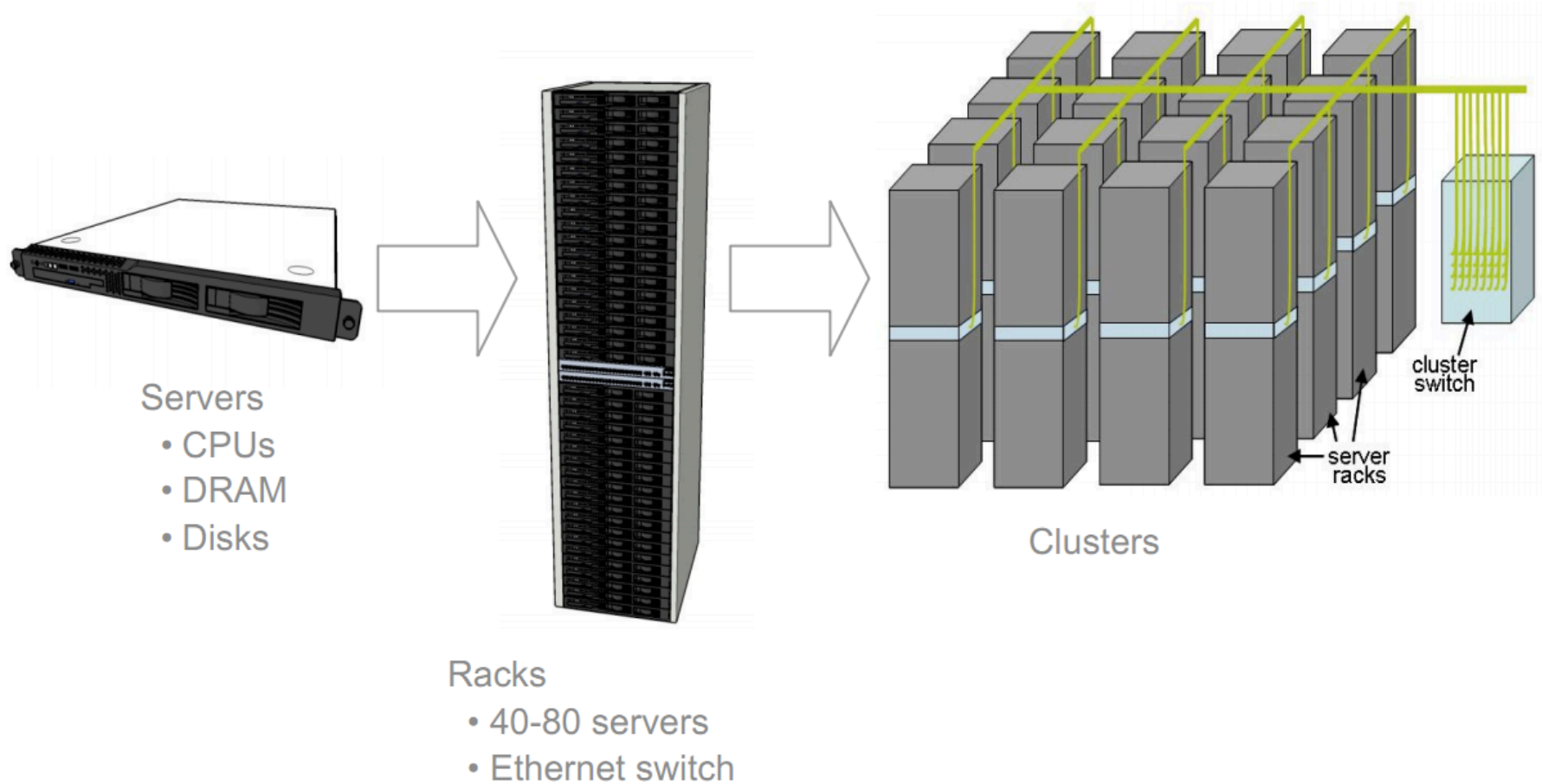
- Service hides software and hardware complexity
- Data engineers can make service scalable to **millions or even billions** of users
- Allows for modular interaction
  - Results from search, maps, Q&A, etc

## Layering in Web Services

- Services widely exposed on the web, accessible via **HTTP**
- **Proxies** route requests to multiple back-end services and join results
- Services themselves can be implemented on **distributed and parallel architectures**



# The Machinery



# The Joys of Real Hardware

Typical first year for a new cluster:

- ~0.5 **overheating** (power down most machines in <5 mins, ~1-2 days to recover)
- ~1 **PDU failure** (~500-1000 machines suddenly disappear, ~6 hours to come back)
- ~1 **rack-move** (plenty of warning, ~500-1000 machines powered down, ~6 hours)
- ~1 **network rewiring** (rolling ~5% of machines down over 2-day span)
- ~20 **rack failures** (40-80 machines instantly disappear, 1-6 hours to get back)
- ~5 **racks go wonky** (40-80 machines see 50% packetloss)
- ~8 **network maintenances** (4 might cause ~30-minute random connectivity losses)
- ~12 **router reloads** (takes out DNS and external vips for a couple minutes)
- ~3 **router failures** (have to immediately pull traffic for an hour)
- ~dozens of minor **30-second blips for dns**
- ~1000 **individual machine failures**
- ~thousands of **hard drive failures**
- **slow disks, bad memory, misconfigured machines, flaky machines, etc.**

Long distance links: **wild dogs, sharks, dead horses, drunken hunters, etc.**



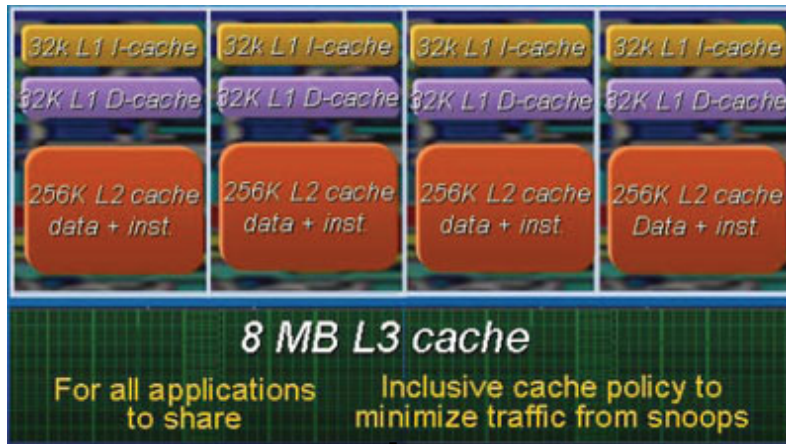
## Reliability & Availability

- Things will crash. Deal with it!
  - Assume you could start with super reliable servers (MTBF of 30 years)
  - Build computing system with 10 thousand of those
  - **Watch one fail per day**
- Fault-tolerant software is inevitable
- Typical yearly flakiness metrics
  - 1-5% of your disk drives will die
  - Servers will crash at least twice (2-4% failure rate)

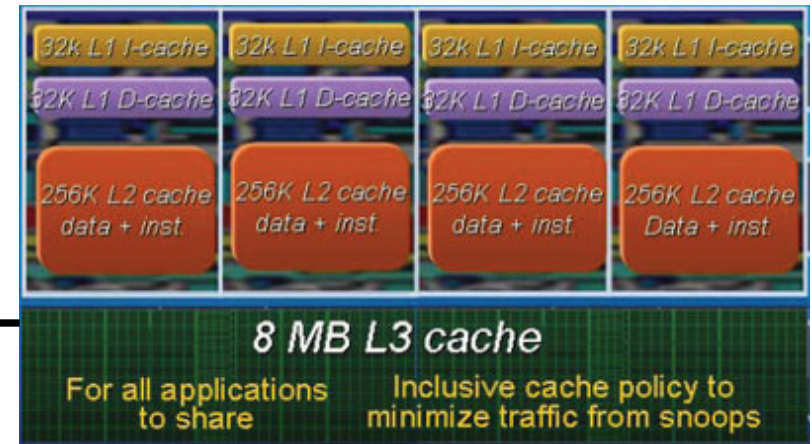


# Complexity lives even inside a single server...

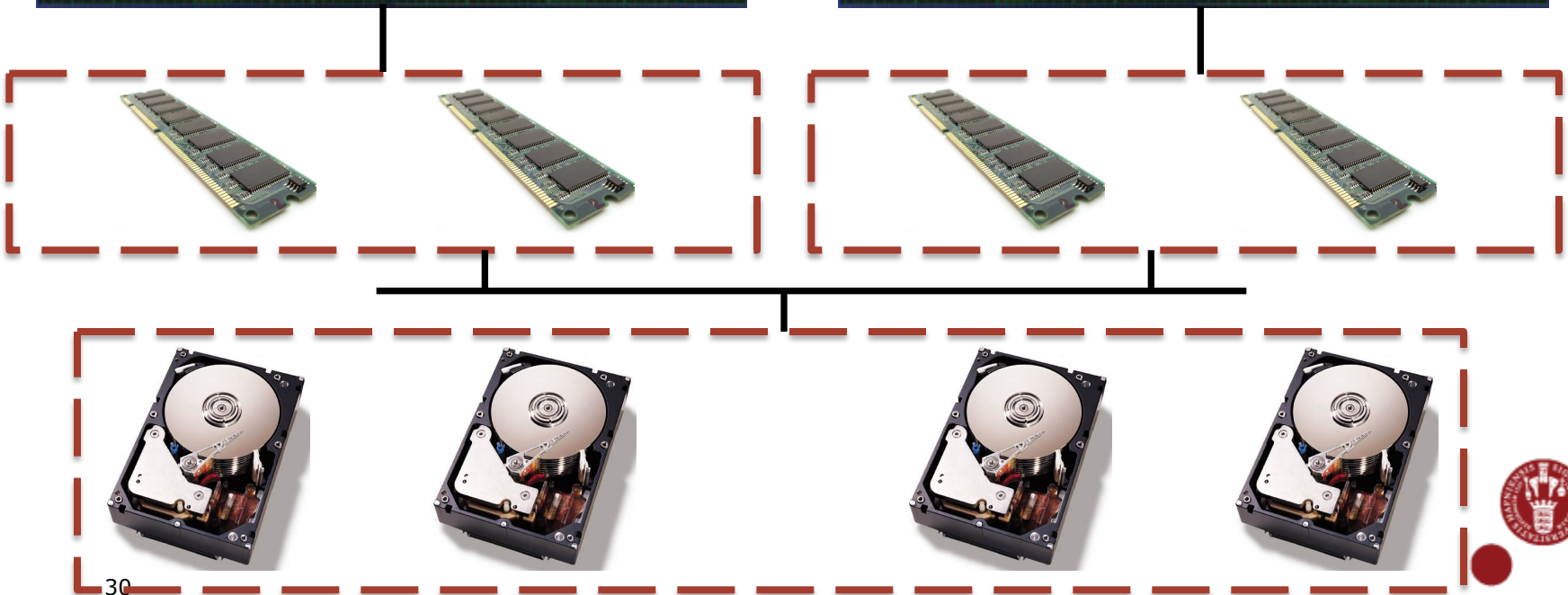
Processor 1



Processor n



...



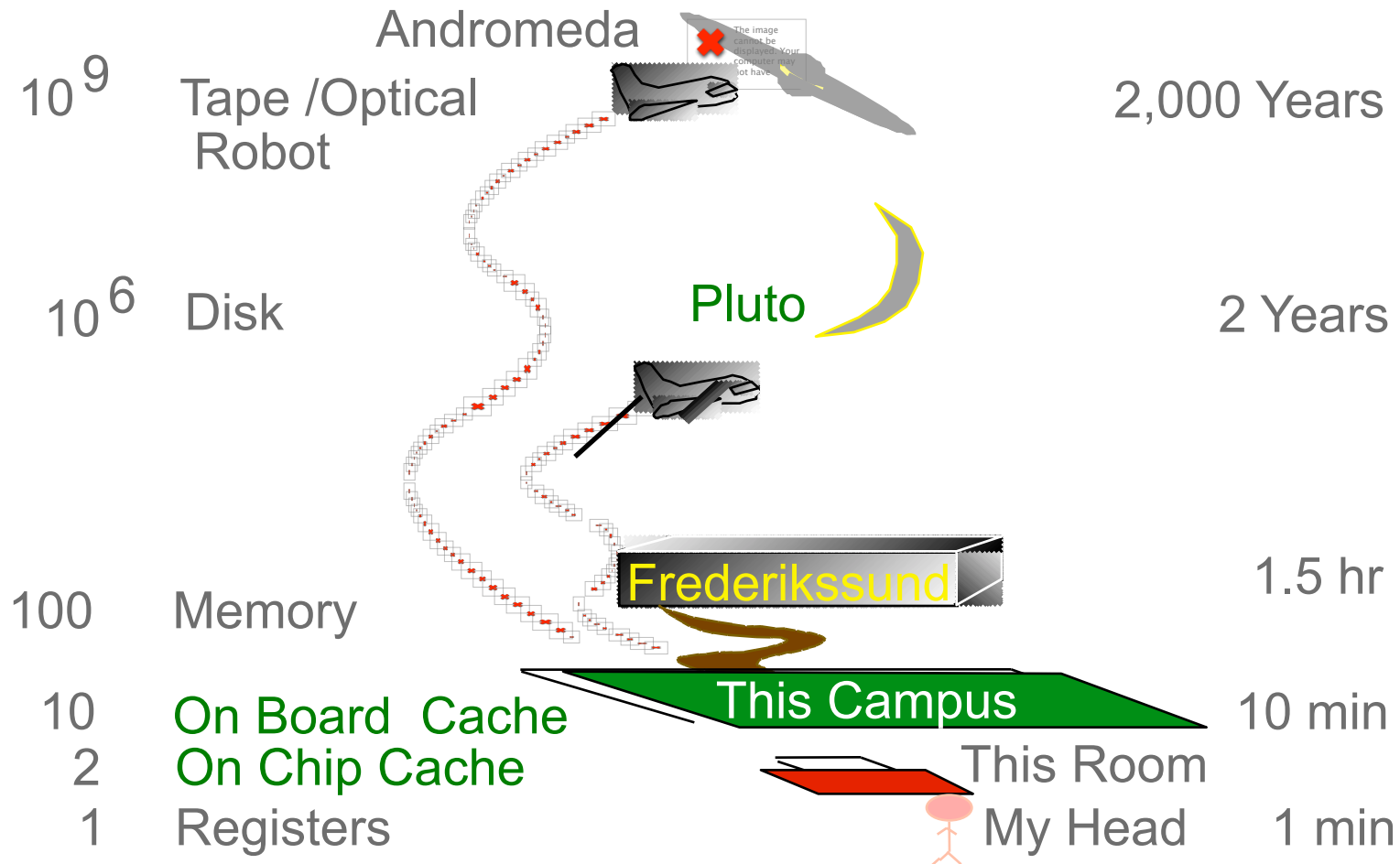
But the picture is not to scale!



What about the size of disk?



# Storage Hierarchy



Source: Gray (partial)



# Common Issues in Designing Services

- **Consistency**
  - How to deal with *updates* from multiple clients?
- **Coherence**
  - How to refresh *caches* while respecting consistency?
- **Scalability**
  - What happens to resource usage if we increase the #clients or the #operations?
- **Fault Tolerance**
  - Under what circumstances will the service be unavailable?



# Research Highlight: Data Platform for Future Cropping Project

- **Traditional Approach to Data Management in Agriculture**

- Build BIG database, e.g., data warehouse
- Lots of time spent mapping schemas, defining what queries to answer
- Inflexible, high cost, limited to specific questions

- **Future Cropping Data Platform**

- Build a service-centric data platform
- Data platform manages and serves geospatial data for analysis services
- Flexible, pay-as-you-go integration of analytic functions
- Separation of concerns: Expertise in scalability for data platform

[https://en.wikipedia.org/wiki/Spatial\\_analysis#/media/File:Snow-cholera-map.jpg](https://en.wikipedia.org/wiki/Spatial_analysis#/media/File:Snow-cholera-map.jpg)

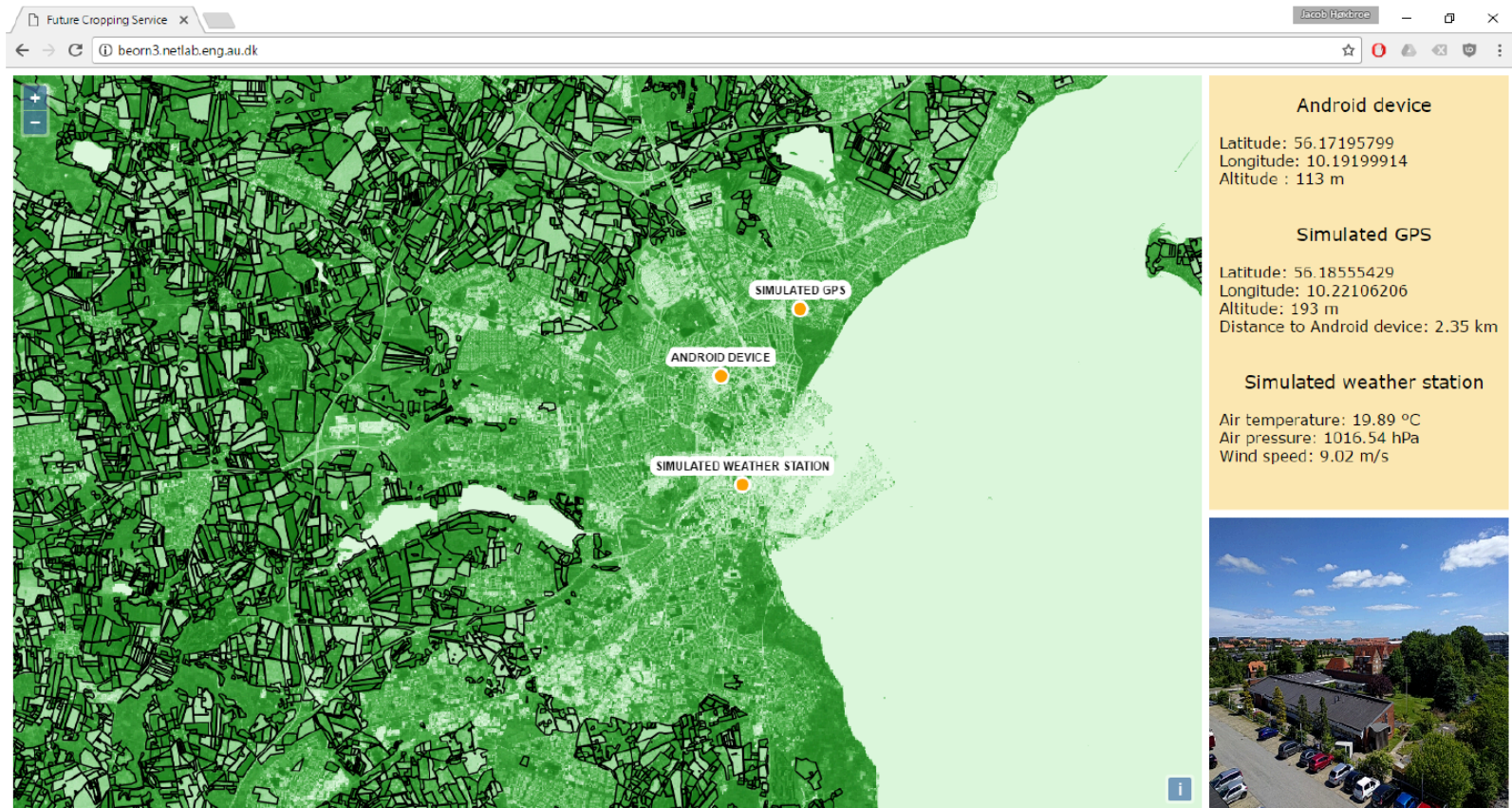


## Spatial Analysis

Work in collaboration with KU PLEN, Aarhus U, and other partners in Future Cropping project; MSc thesis of Mads Engesgaard Jacobsen and ongoing PhD of Yiwen Wang

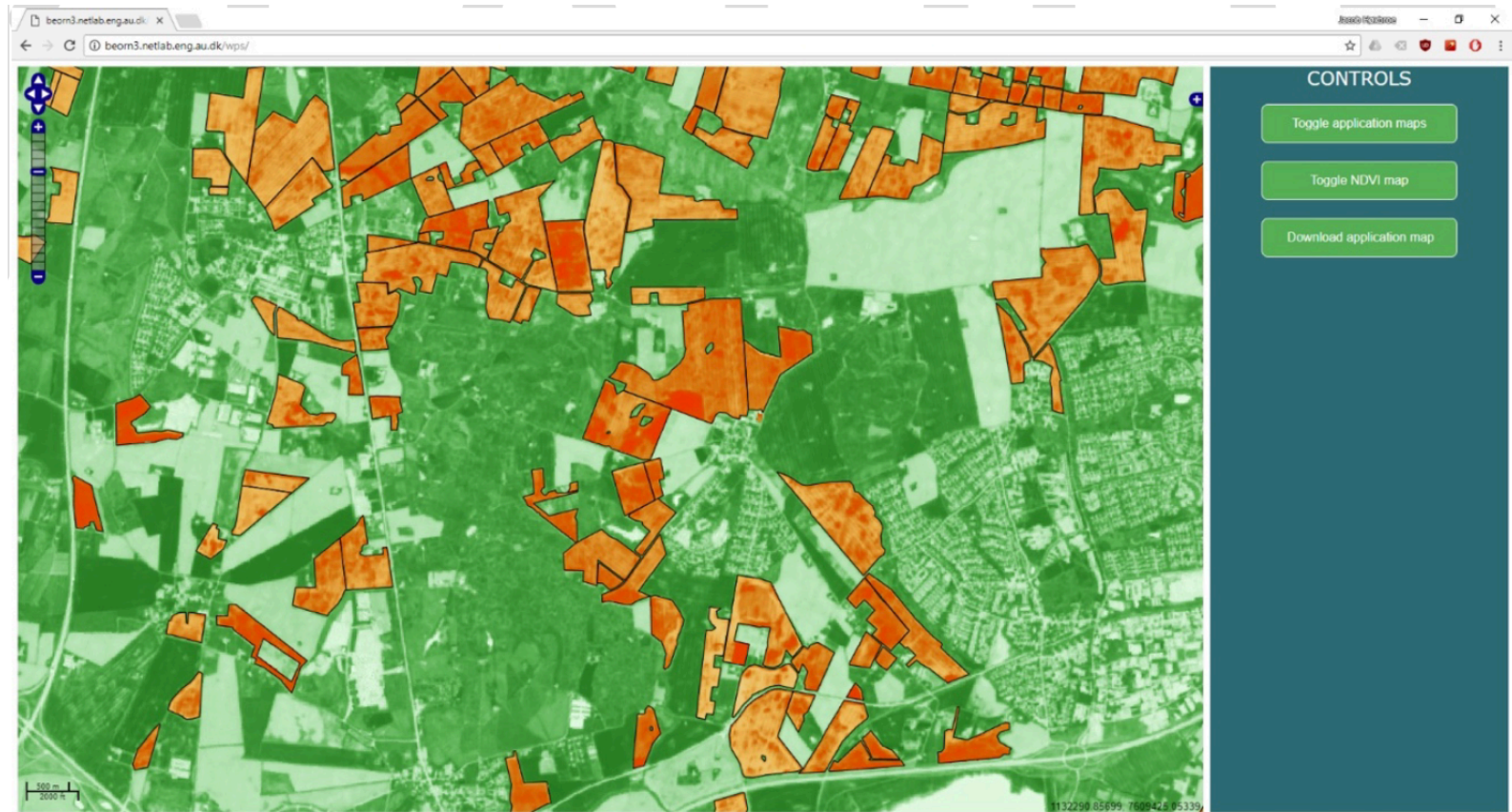


# Examples of Services



- Online streaming data from moving objects overlaid with Sentinel-2 satellite data and field polygons
  - Demo available at: <http://beorn3.netlab.eng.au.dk/>

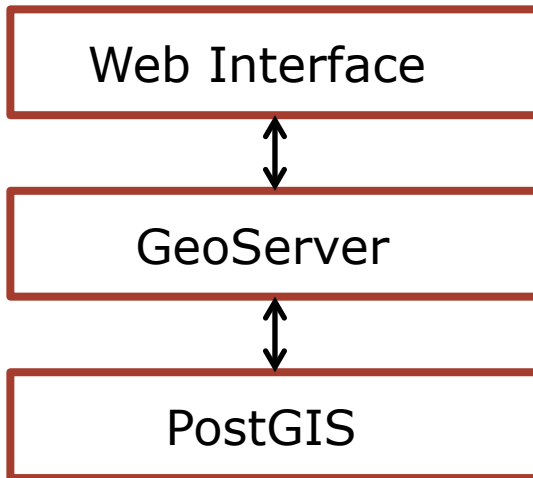
# Examples of Services



- Field polygons and NDVI map with option to download fertilizer application maps
  - Demo available at:  
<http://beorn3.netlab.eng.au.dk/wps/>

# Data Platform Foundations and Trends

- **GeoNode**



- **Trends**

- More users
- More datasets
- More frequent updates
- More analytical services

- **Existing systems not enough for the future**

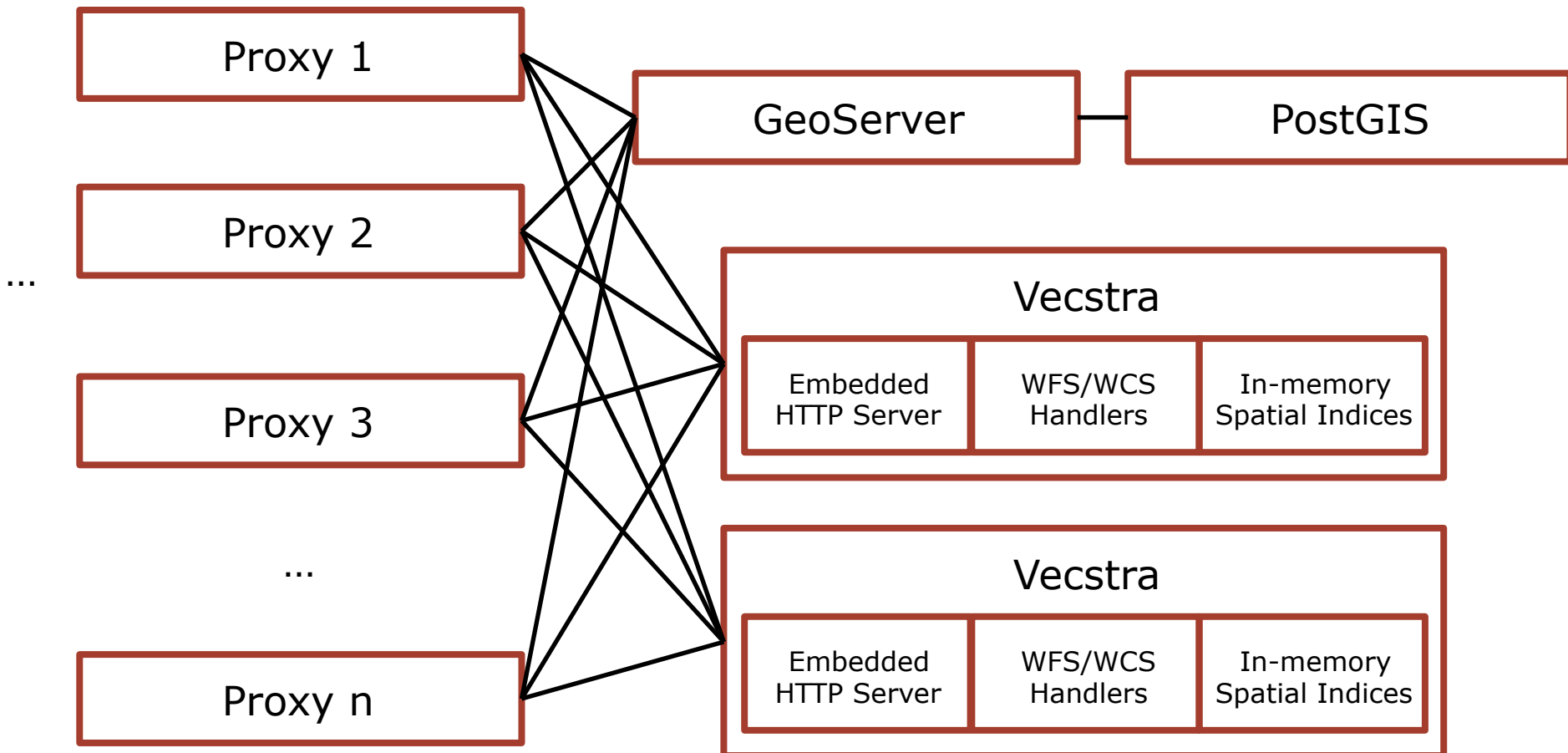
- General caches, e.g., Redis, Memcached, not specialized for geospatial data
- Geospatial caches, e.g., GeoWebCache, only cover subset of protocols (WMS)

## Vecstra: Core ideas

- **Support for *transparent scalability* on concurrent requests**
  - Communicate with cache through subset of WFS/WCS protocols so as to make solution “drop-in”
  - Reverse proxy layer routes to multiple caches or falls back to GeoNode for advanced functionality
- ***Low latency* in serving data**
  - Employ state-of-the-art in-memory spatial indices
  - Revisit algorithms to speed up specific operations, e.g., geometry intersection, counting queries
- ***Efficiency* in use of computational resources**
  - Design cache multi-threading for performance on multi-core servers



# Vecstra Architecture



## Initial Evaluation: WFS

- **Workload modeling Future Cropping**

- *Crop Status Service*

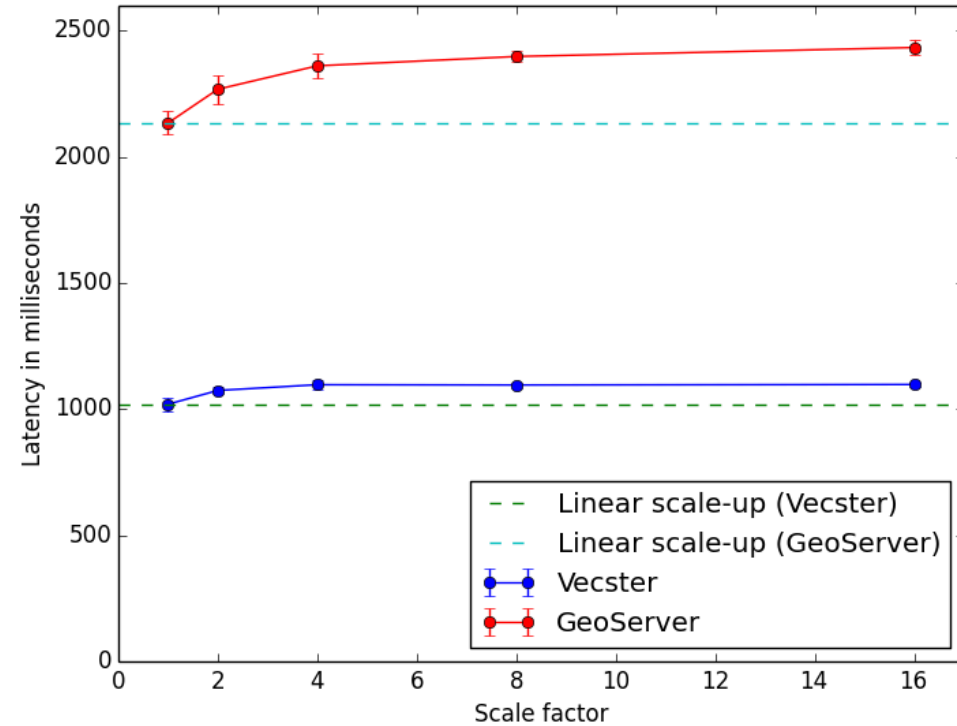
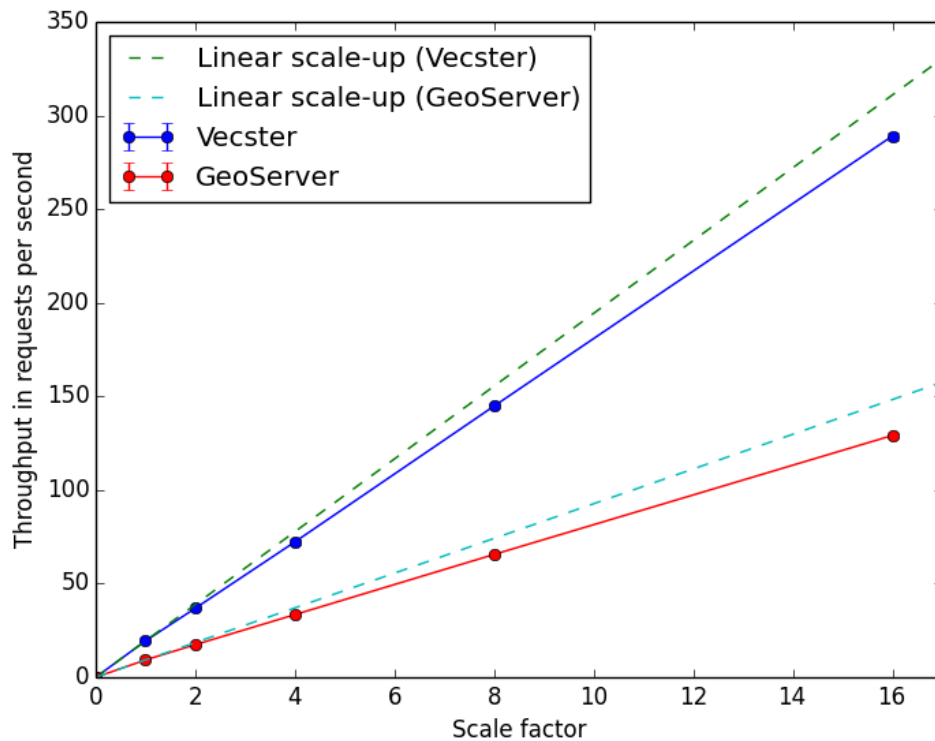
- Field-sized query windows scaled between 0.5x-2x
    - Fields picked uniformly at random, window centered in point within bounding box of field
    - Layers: field polygons, topography, soil, rain distribution (future: also NDVI, climate)
  - No updates for now; need to model update patterns

- **Single-node multi-core server**

- Vector and NDVI layers in data platform as of early 2017 take roughly 23 GB → *in-memory processing*
  - *Server: 16 cores; 2 sockets; HT not used; 128 GB RAM*
  - Thread affinity or taskset used to limit cores used
  - 20 client threads per server core in separate machine; 10Gbit Ethernet

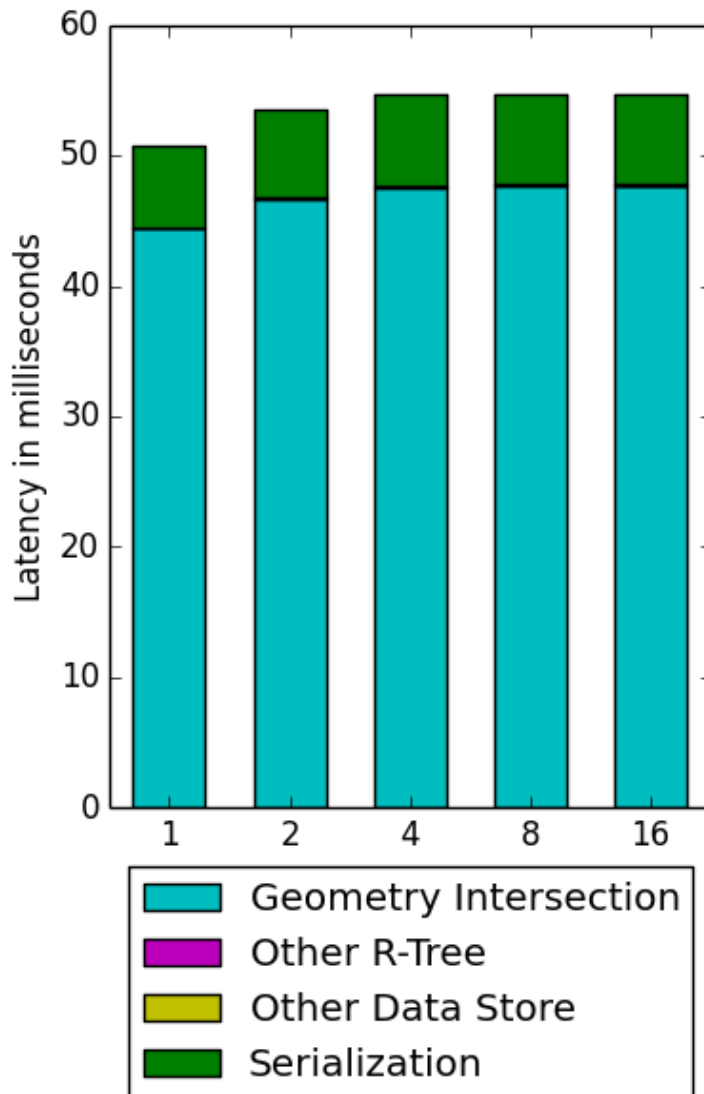
# Initial Evaluation: WFS Spatial Range Query

## Multi-Core Scalability



- Both systems scale relatively well with increasing numbers of cores in high-load, read-only setup
- Under minimal tuning of both systems, Vecstra shows promise of delivering better latency and scaling efficiency

# Initial Evaluation of WFS Spatial Range Query: Where does the time go?



- Server-side latency (HTTP not included)
- Recall we increase client threads (20) to push utilization up
- On the server side, most of the query time is spent in geometry intersection operations, i.e., in refine step of filter-refine scheme
- More work to be done here 😊

# Lesson 4: Telemetry Turns Behavior into Data



## Usage logs are data!

- **Web log records queries of a web service**
  - User access patterns include **spatial and temporal information**
- **Model user attention and skew in access patterns**
  - For better caching, deployment of computational infrastructure
  - For detection of patterns leading to bias in business decisions
    - Products that users most looked at
    - Breeds that users pay “too much” attention to



Steve Jurvetson - <https://www.flickr.com/photos/jurvetson/162116759>

## Research Highlight: TileHeat

- **Case study of The Digital Map Supply**
  - Most popular WMS web service of Danish National Survey and Cadastre, now Danish Geodata Agency
- **Issues**
  - Render service is slow to compute tiles (for some map services)
  - Bulk data updates by Danish municipalities
  - In combination: Bad performance (for some map services)
- **Data we have analyzed**
  - Request log last 5 years: ~1B requests total
  - Q4 2011 log for most popular map service: ~800K requests per day



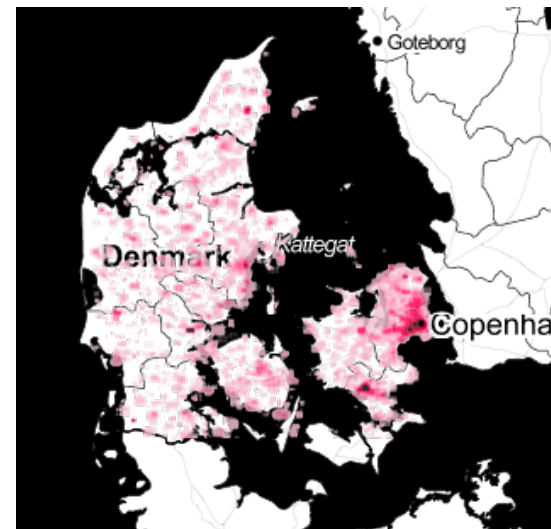
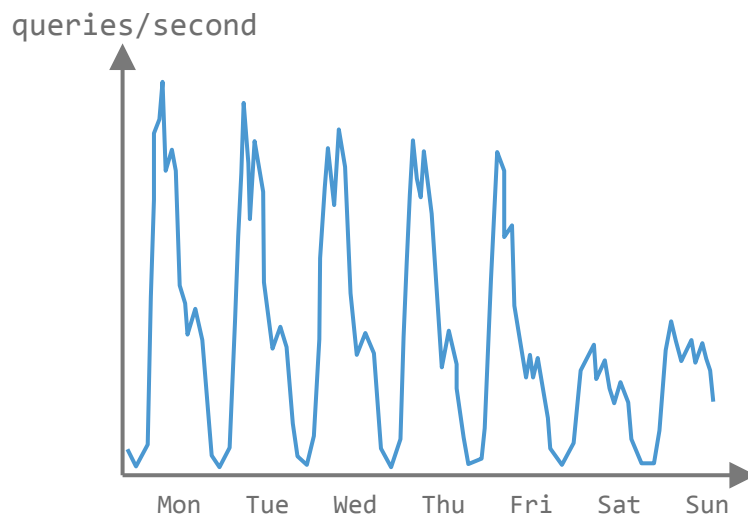
Geodatastyrelsen

<http://gst.dk/>

Work done in collaboration with P. K. Kefaloukos and M. Zachariasen, results in ACM SIGSPATIAL GIS 2012



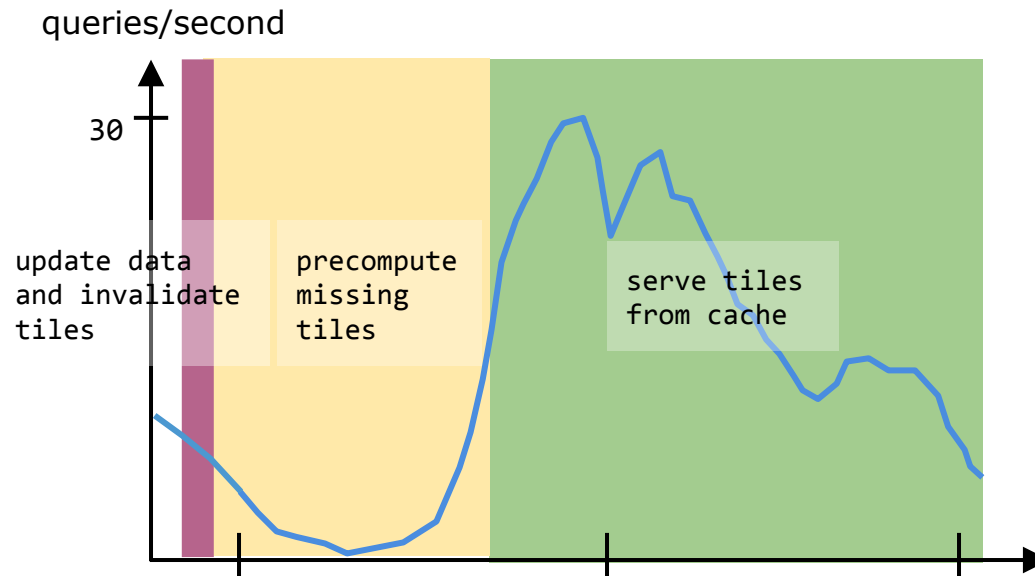
# Exploitable Properties



Prediction of where people will look on the map

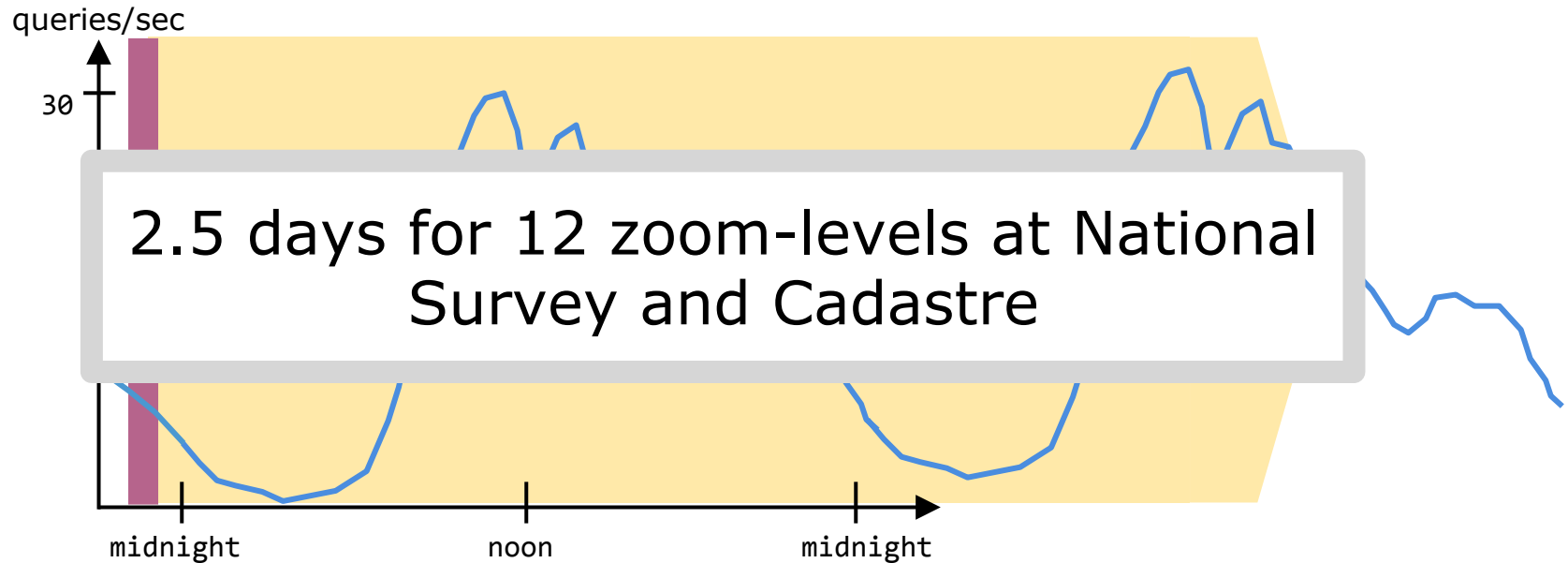
- 1) Seasonal variation in load (24-hour, week)
  - 2) We can predict\* the tiles people will tend to request
  - 3) Strong skew in requested tiles
- \* ) For the maps we have studied, but not necessarily a priori!

# Ideal Situation



- Schedule massive data updates during low load
- Time to refresh cache with new tiles before peak load
- Serve tiles from cache during peak load

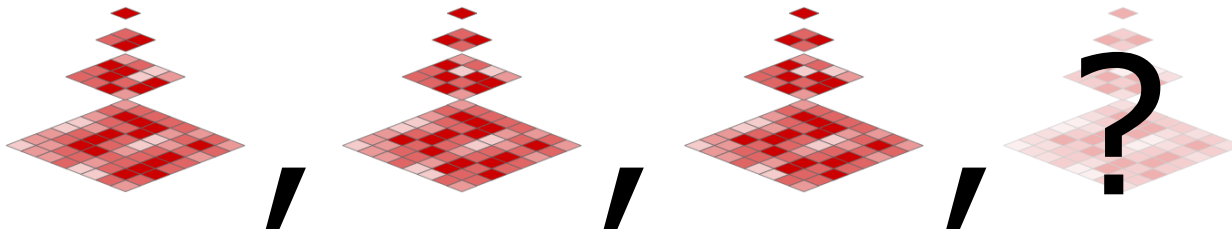
## Problem



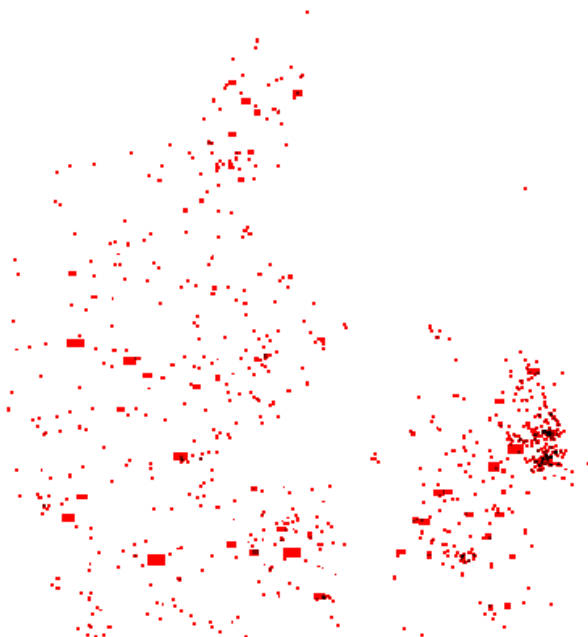
- It takes a long time to generate all tiles
- $O(4^m)$  for  $m$  zoom-levels

## Our goal

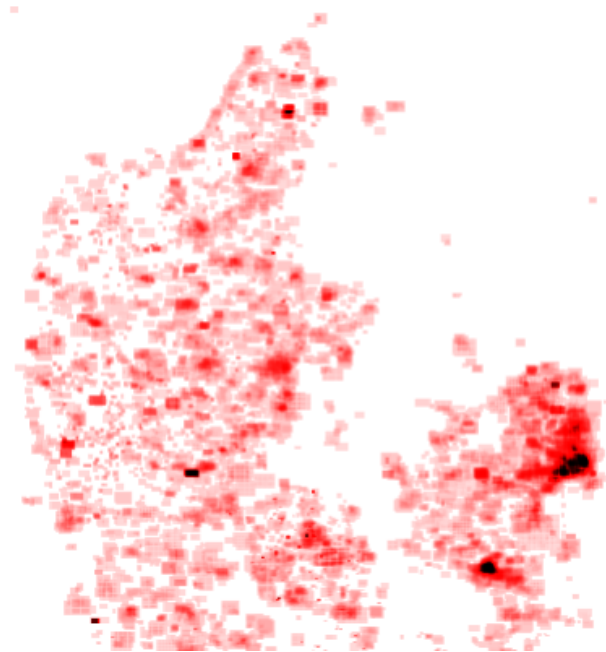
- Predict the heatmap of tomorrow



# Heat Dissipation: A Real Example

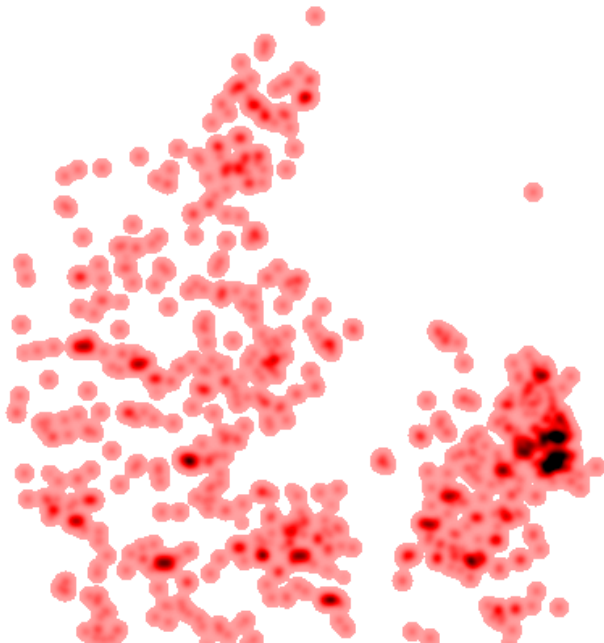


Heatmap based on a small sample of requests (real data)

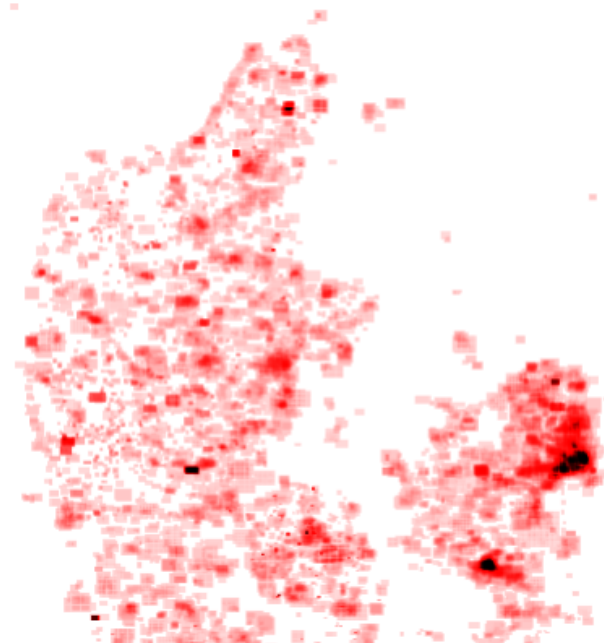


Real heatmap

# Heat Dissipation: A Real Example



After dissipation



Real heatmap

- In addition to heat dissipation, exponential smoothing model used in predictions to capture variations in time

A map of Scandinavia with three green rectangular boxes highlighting specific regions. The top box is over Sweden, with the label 'Göteborg' to its right. The middle box is over Denmark, with the label 'Denmark' centered within it. The bottom box is over the eastern part of Denmark, with the label 'København' to its right. The label 'Copenhagen' is also visible below the bottom box.

# Lesson 5: Embed Intelligence in Services



From querying and calculation to prediction

- **Types and examples of**

**BEWARE: Data  
quality is the  
primary success  
factor!**

obtained so far for all farms  
in the country together with  
my planned interventions for  
this year



)#/media/  
Wykoff.JPG

# Research Highlight: GANDALF project

- **Soil Contamination Management in Urban Areas**

- Limited sets of pre-selected indicators of potential pollutants in chemical analyses

- **GANDALF**

- Leverage historical data in spatial interpolation model
- Enrich existing techniques with machine learning approaches, make data-driven decisions, e.g., for where to sample next
- Move towards untargeted chemical fingerprinting with high dimensionality and merging with historical data

[https://en.wikipedia.org/wiki/Spatial\\_analysis#/media/File:Snow-cholera-map.jpg](https://en.wikipedia.org/wiki/Spatial_analysis#/media/File:Snow-cholera-map.jpg)



## Spatial Analysis

Work in collaboration with KU PLEN, MOE, KMC Nordhavn, and other partners in GANDALF project

Lots of machine learning work going on at DIKU!

## Lessons From Managing Geospatial Data

- **Challenge: Big variation in data formats and volume**
  - Lesson 1: “Cheap” vs. “expensive” data
  - Lesson 2: The rise of standardization, open-source software, and large geospatial datasets
- **Challenge: Large amount of users and potentially complex simultaneous requests**
  - Lesson 3: From software to services
  - Lesson 4: Telemetry turns behavior into data
- **Challenge: Much labor needed to derive knowledge from varied data**
  - Lesson 5: Embed intelligence in services



## Conclusion

- **Spatial Applications & Challenges**
- **From Challenges to Lessons**
  - Lesson 1: “Cheap” vs. “expensive” data
  - Lesson 2: The rise of standardization, open-source software, and large geospatial datasets
  - Lesson 3: From software to services
  - Lesson 4: Telemetry turns behavior into data
  - Lesson 5: Embedded intelligence in services
- **Workshop**
  - Groups take lessons as input and discuss how they can be applied to plant phenotyping area
  - Groups summarize discussion work and present in plenum

Thank you!



# Background Information



## About the Speaker: **Marcos Vaz Salles**

- Associate Professor, University of Copenhagen (**DIKU**)
  - Postdoc: Cornell University
  - PhD: ETH Zurich
- **Expertise:** Database Systems
  - In-memory databases
  - Spatial data
  - Information Integration
  - Cloud Computing
- Co-leader of **Data Management Systems (DMS) Lab**
- **Ongoing Collaborations**
  - Future Cropping consortium: precision agriculture
  - GANDALF consortium: environmental management
  - IDAS: Industrial Data Analysis Service
  - HIPERFIT center: financial apps, risk management

