# Visualizing health data – from fundamental research to successful applications
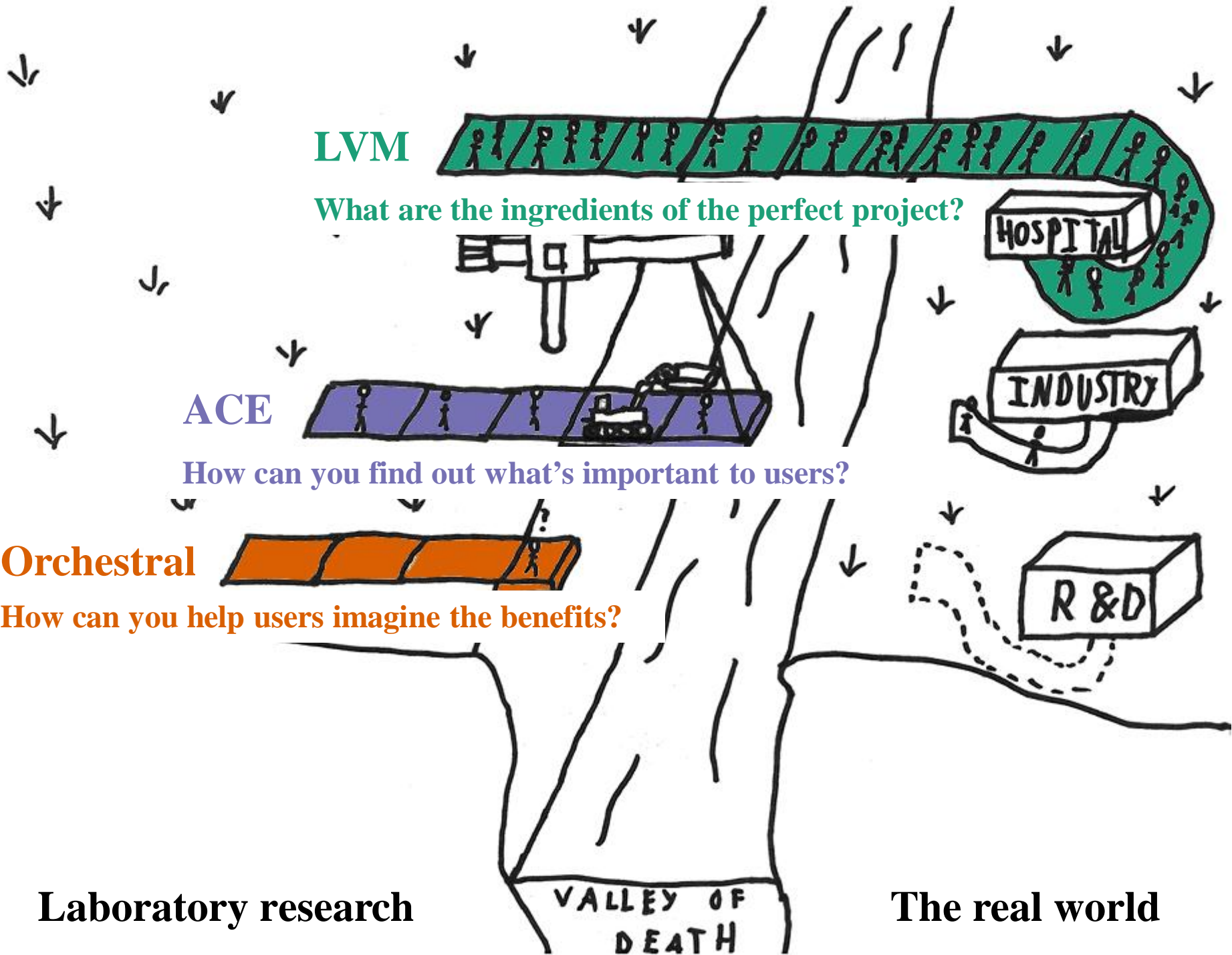
**Professor Roy Ruddle**

**University of Leeds & Alan Turing Institute**

**https://raruddle.wordpress.com/**

The Alan Turing Institute

UNIVERSITY OF LEEDS

LVM

What are the ingredients of the perfect project?

ACE

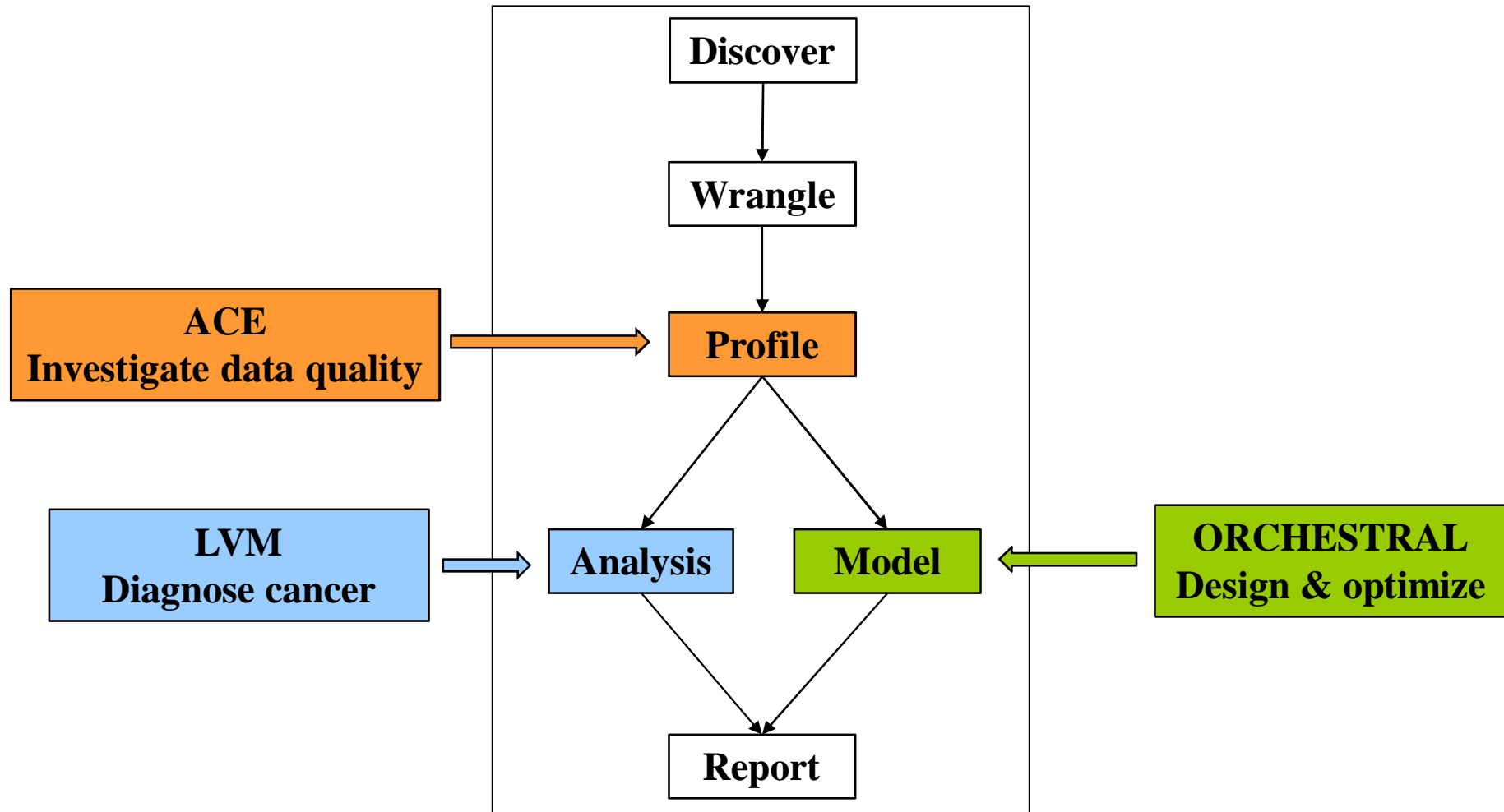How can you find out what's important to users?

Orchestral

How can you help users imagine the benefits?

Laboratory research

The real world

VALLEY OF DEATH

# Data science workflow[1]



[1]Alspaugh et al. (2018). *IEEE TVCG*.

# ACE
## A novel tool for investigating missing data

# Big picture

- **Data wrangling & profiling take 50 – 80% of data scientists' time**
- **Many tools for investigating data quality**
  - **But they don't meet users' requirements**
- **Users lack of knowledge & rigour in data cleaning[1]**
- **Visualization methods for data quality**
  - **Limited research[2]**
  - **Unrealistic evaluation (toy datasets)**

---

[1]Visualizing the quality of data https://tinyurl.com/VizDataQuality

[2]Arbesser, et al. (2017). *IEEE TVCG*; Gotz, & Stavropoulos. (2014). *IEEE TVCG*; Gratzl, et al. (2013). *IEEE TVCG*; Gschwandtner, et al. (2014). *Proc. I-KNOW*; Kandel, et al. (2012)). *Proc. AVI*; Noselli, et al. (2017). *Proc. HEALTHINF*; Ruddle & Hall. (2019). *Proc. HEALTHINF*; Tennekes, et al. (2011). *Proc. NTTS*; Unwin, et al. 1996). *Comp. & Graph. Stat*; Xie, et al. (2006). *Proc. IEEE VAST*; Zhang, et al. (2014). Information Visualization

# How can you find out users' requirements?

- **Tensions in applied research**
  - **Useful tool vs. novel research**
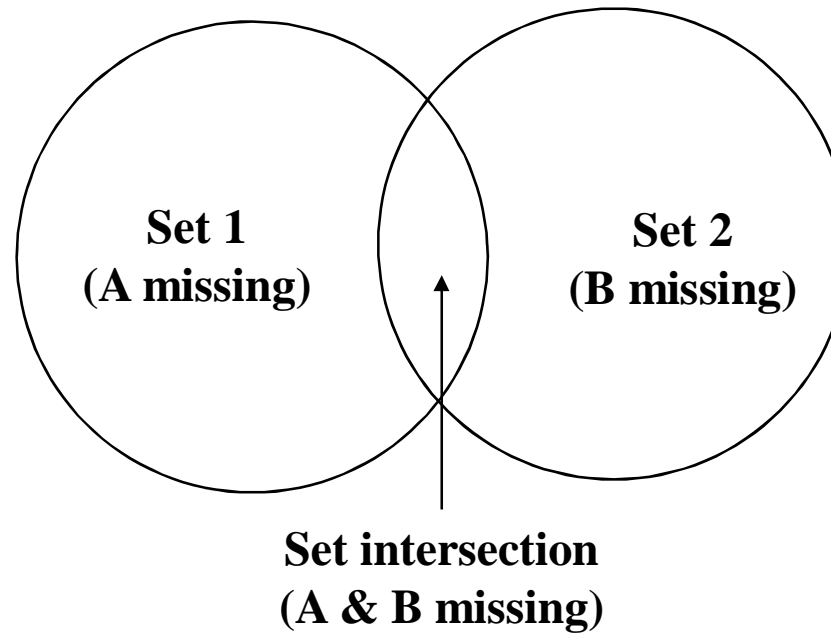  - **Market research vs. requirements analysis**

| Three steps | Methods |
|---|---|
| Find out current situation | Questionnaire, interview, documentation, example data:<br>• What analysis steps are involved?<br>• What do you already do well?<br>• What do you know? |
| Explore what's needed | Ask what would you like to do, but cannot do today?<br>• What is hard or time-consuming (barriers & bottlenecks)?<br>• What assumptions/simplifications are you forced to make?<br>• Why don't current analysis tools solve these difficulties?<br>• Let your self dream … |
| Check your understanding | Workshop<br>• Encourage corrections & comments<br>• Propose solutions (storyboard; throw-away prototype) |

# NHS Digital

- **Provides information, data and IT systems for National Health Service in England (£400 million)**

- **Current situation**
  - **Collect patient-level data from every NHS hospital**
    - **E.g., Admitted Patient Care (APC) data**
      - **500 fields and 20 million records/year**
  - **Mature data cleaning process, including**
    - **Business rules for data correction & validation**
    - **Threshold for missing values (only ≈8 fields)**
    - **Feedback to hospitals**

- **What's needed**
  - **Explore data quality patterns involving multiple fields**
  - **Exclude expected patterns, to reveal the unexpected**
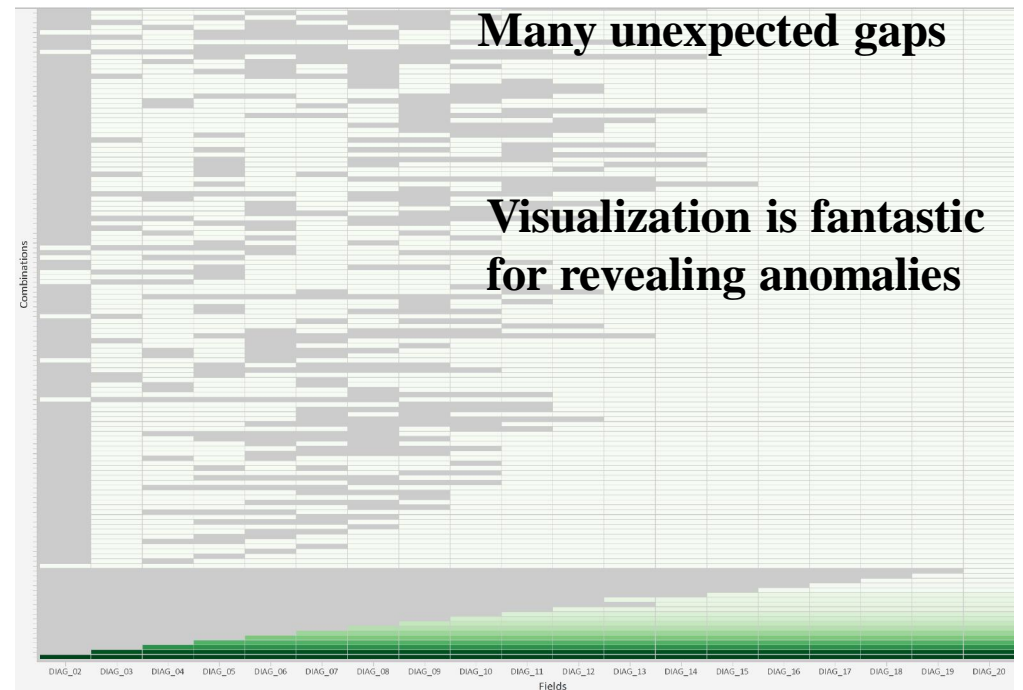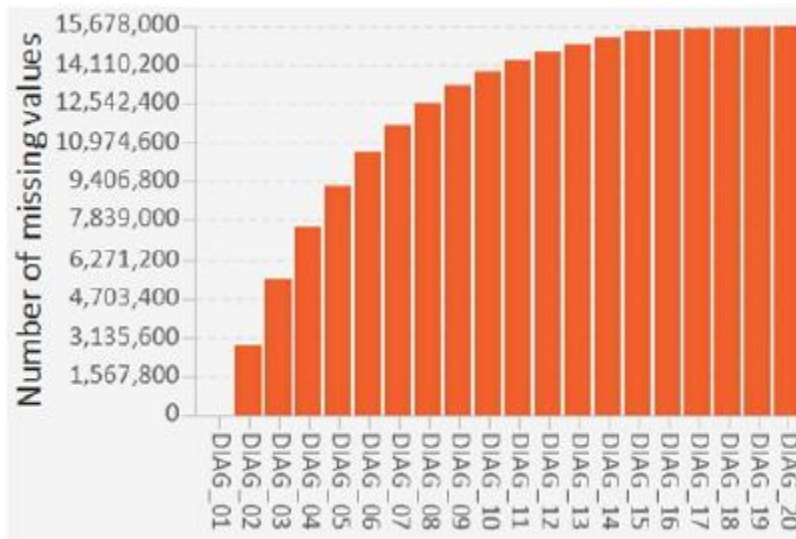  - **Develop new business rules**

# Novel set visualization tool

| Field A | Field B | Field C |
|---------|---------|---------|
| 101     |         | M       |
| 102     |         | F       |
|         |         | M       |
|         | 99      | F       |
|         | 68      | M       |

Set 1
(A missing)

Set 2
(B missing)

Set intersection
(A & B missing)

- **Scalable design**
  - **20 million records**
  - **500 fields**
  - **500,000 combinations of missing values**
- **Achieved using well-known techniques**
  - **Bar charts, heat maps and histograms**
  - **Reduce learning curve (avoid unnecessary novelty)**

# Admitted Patient Care (APC) example

- **20 fields for diagnostic codes**
  - Missing more often from DIAG_01 to DIAG_20

# Actionable insights

- **Widespread implications for data cleaning rules**
- **Gaps in diagnostic codes**
  - **Only 2000 records**
  - **85+ % from one admission method in specific hospital**
    - **Improve data quality via established mechanism**
- **Gaps in operation codes**
  - **2500 records**
  - **May affect NHS Payment by Results system for hospitals**
- **Millions of missing dates**

**Laboratory research**

**The real world**

NHS Digital &
6 other end-user
partners

ACE

Diagnostic codes

Workflow
Too busy

HOSPITAL

INDUSTRY

R & D

VALLEY OF DEATH

# Orchestral
## Visualizing genomics algorithms

# Big picture - genomics

- **Mature tools and pioneer in "big data"**
- **But unimaginative visualization**
  - Massive over-plotting of data points
  - Have to pan thousands of times
- **Example application: Breast cancer**
  - Clear (subjective) differences between cohorts
  - Need to understand differences to design statistical test

# Hypothesis

- **Large high-resolution displays could transform scientists' ability to find patterns in genomic data**



54 megapixel Powerwall (3 x 1.3 metres)

# How can you help users imagine the benefits?

- **Throw-away prototypes**
  - **Giant image (get user "pull")**
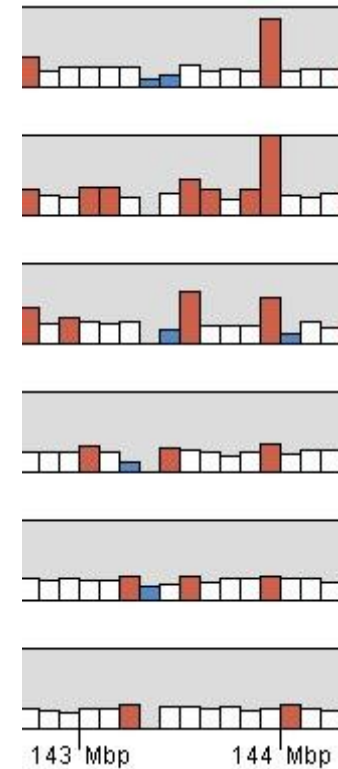  - **Static visualization (spatial compression was too radical)**
  - **Interactive proof of concept**

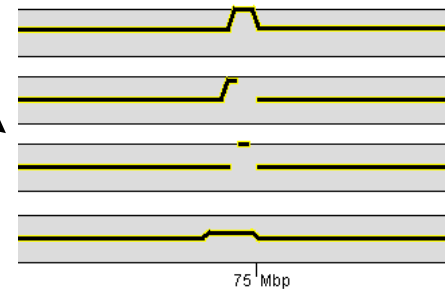# Current visualization vs. Orchestral



**Current visualization**

Copy number variation
(100 kilobase windows)

Segmented data
(noise removed for
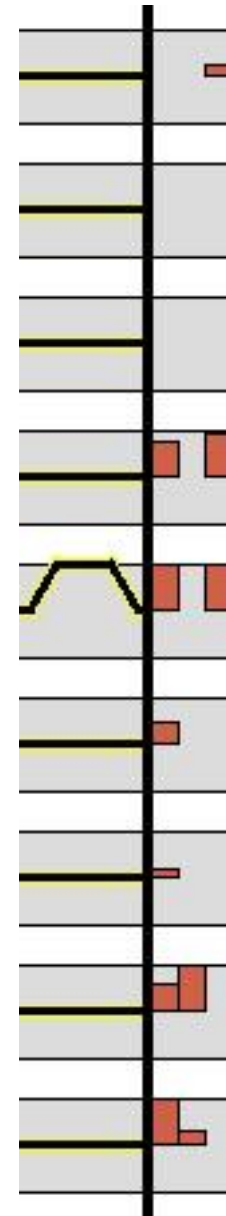statistical analysis)

**Orchestral**

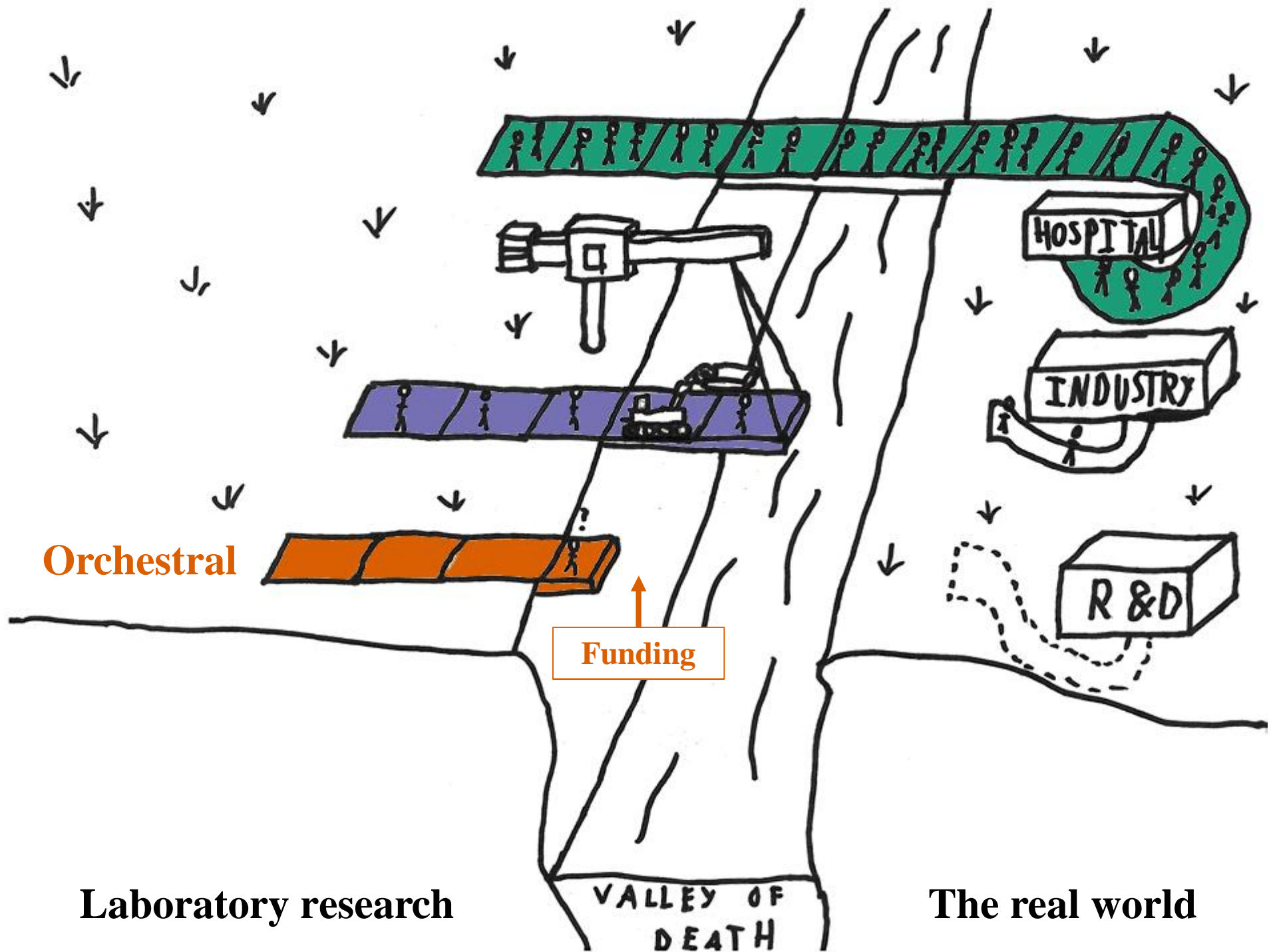# Open the black box by visualizing detail in context

- **"Data looks abnormally similar, almost identical"**
  - Processing error (incompatible steps)[1]
- **Smoothing algorithm removes common feature**

24 megapixel workstation

[1]Ruddle et al. (2013). *Proc. Biovis.*

Orchestral

Funding

Laboratory research
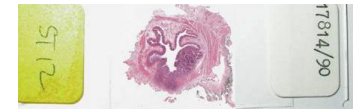
VALLEY OF DEATH

The real world

HOSPITAL

INDUSTRY

R & D

# Leeds Virtual Microscope

## Diagnosing cancer from Amazon-sized images

# What is pathology?

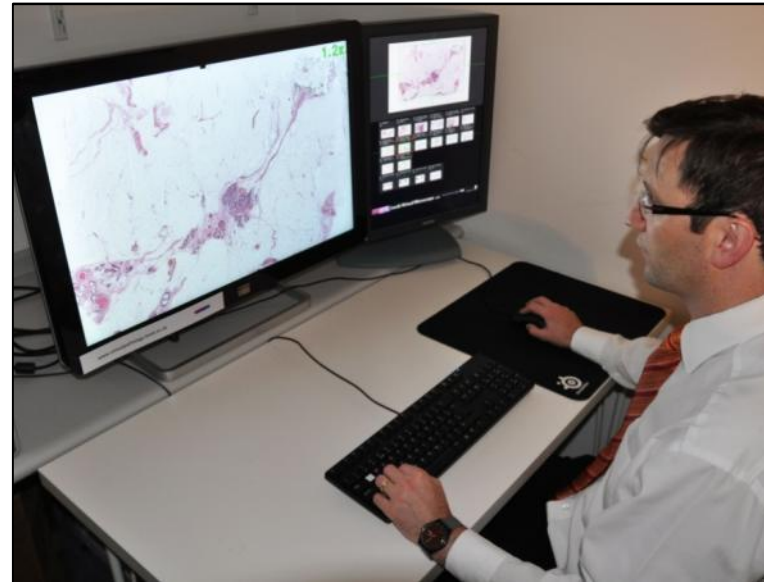"**Pathologists diagnose cancer by using a microscope to examine glass slides that contain thin sections of human tissue**"

- **Large-scale operation**
  - **40 consultants (Leeds Teaching Hospitals)**
  - **150,000 slides/year, at 25 – 400× magnification**
- **The slides can be digitised for viewing on a computer**
  - **Advantageous for 2nd opinions, long-term survival and computer-assisted diagnosis**
  - **But it takes pathologist 60% longer to make a diagnosis[1]**
  - **Each slide is enormous**
    - **10 gigapixels (an "Amazon" of image data)**
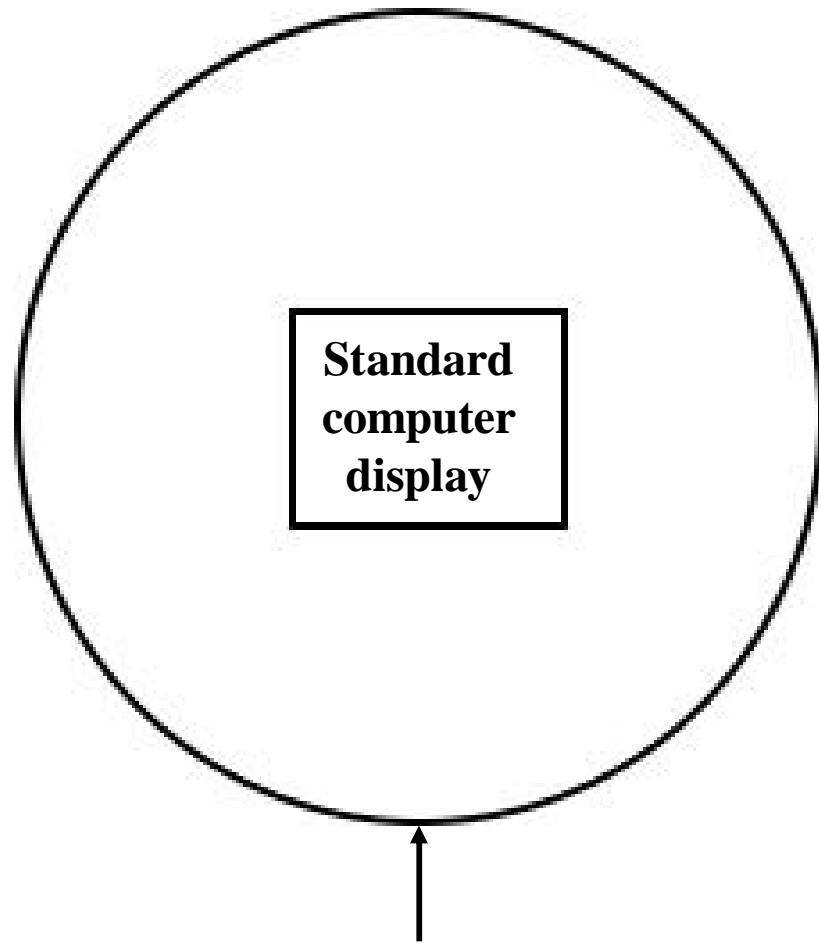
- **http://www.youtube.com/watch?v=oZGkhKkDG5o**

[1]Treanor & Quirke. *Pathological Society Glasgow, July 2007.*

# Why is diagnosis 60% slower?

- **Three reasons**
  - Standard computer displays are too small
  - User interfaces of commercial products are inefficient
  - Doctors lack experience & training with digital slides
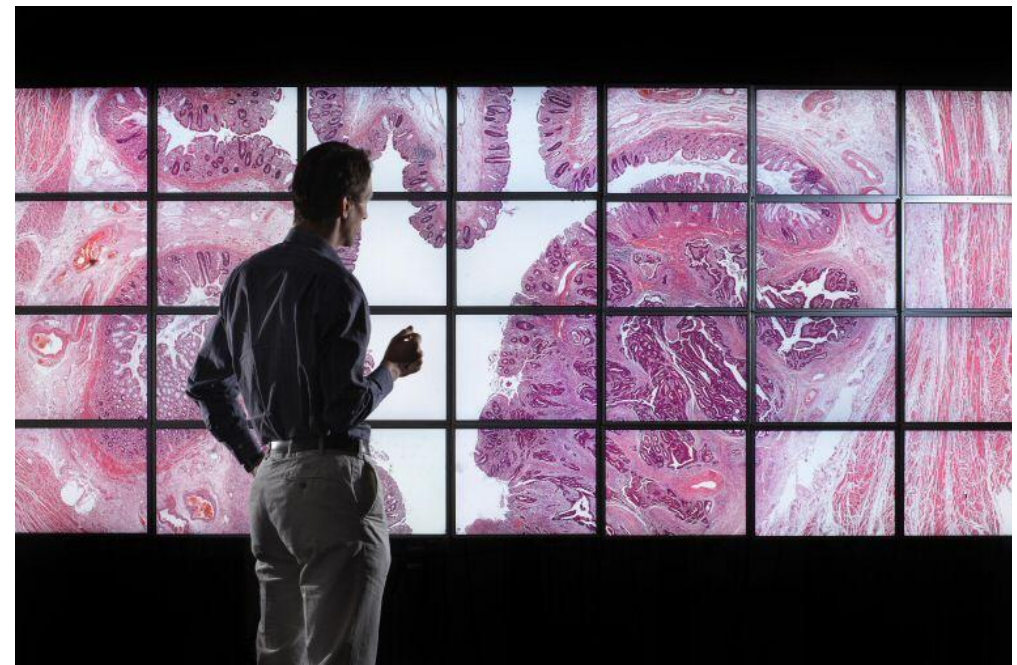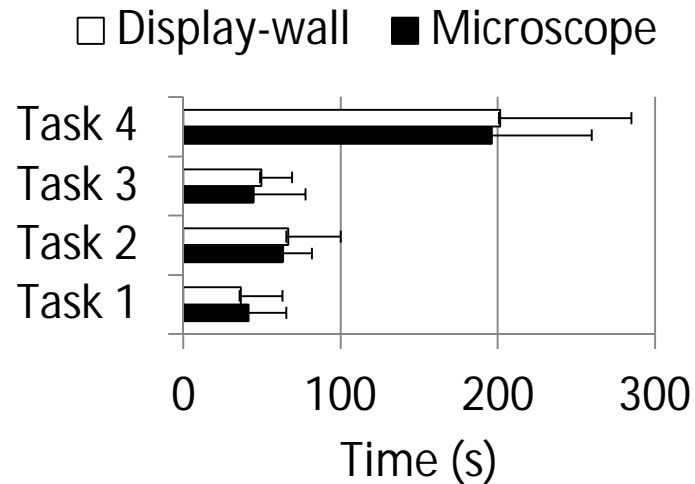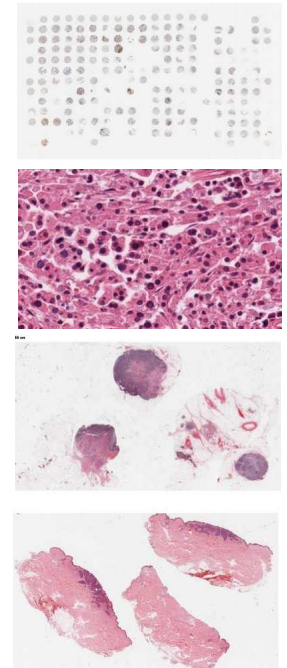
# Standard displays – like looking through a keyhole



Standard computer display

↑
Microscope field of view (48°)

# Solution: 54 megapixel Powerwall[1]

- **6× large field than a microscope, with**
  - 3200 x 2400 pixel "thumbnail"
  - Gamepad user interface
- **Microscope vs. Powerwall evaluation**
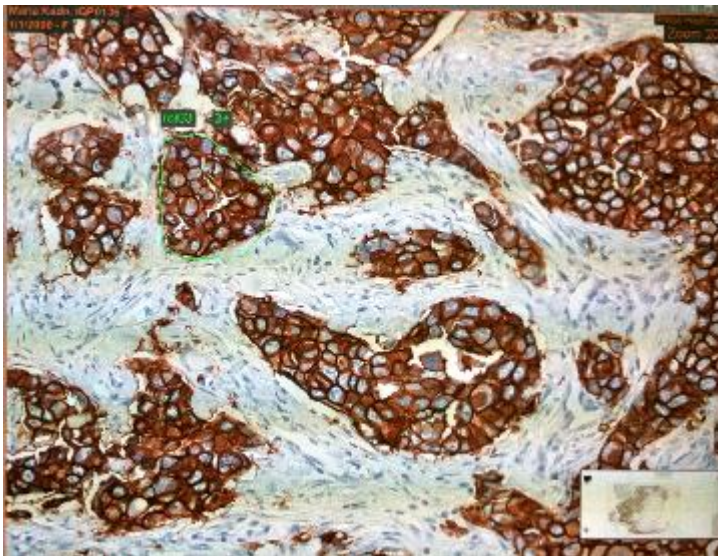  - 4 consultants & 4 trainees
  - Only a few minutes of training



□ Display-wall   ■ Microscope



[1]Treanor et al. (2009). *Histopathology.*

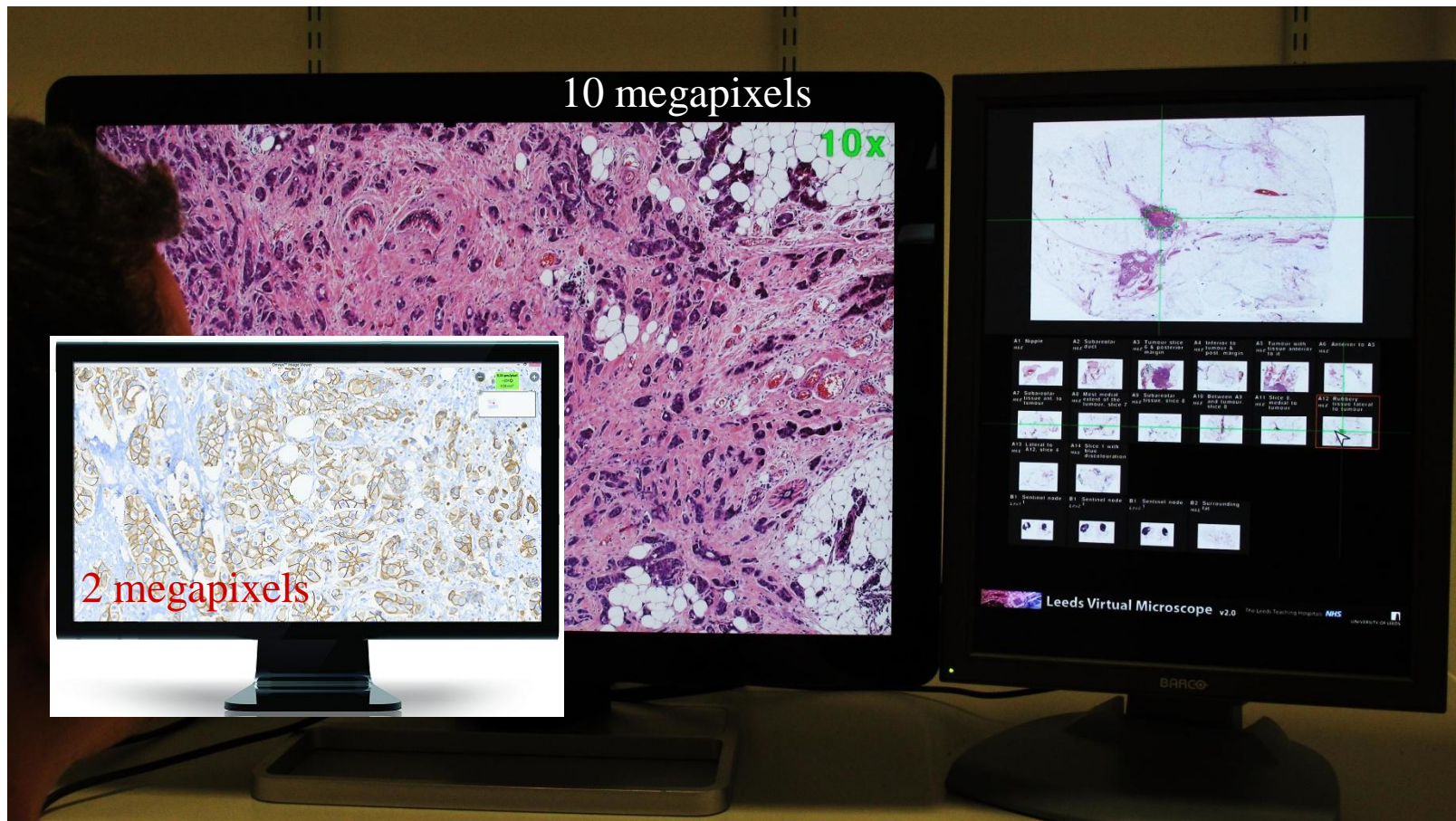# Existing user interfaces – based on Photoshop/Google Maps

- **Glass-sized thumbnail**

- **Real-time interaction**
  - But thousands of panning movements

- **Thmbnail scale difference**
  - 1 : 1200 pathology
  - 1 : 30 Google Maps (established guidelines[1])
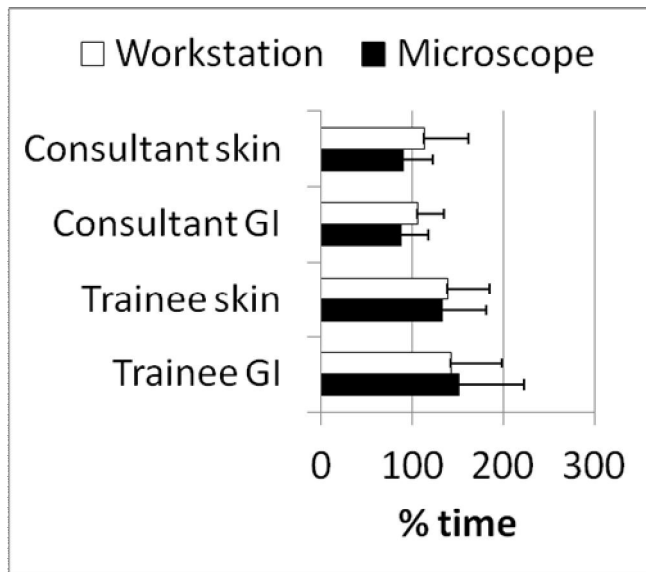




[1]Shneiderman (1998). *Designing the user interface.*

# LVM solution

- **10 megapixel medical-grade display**
- **One third of space devoted to overviews**
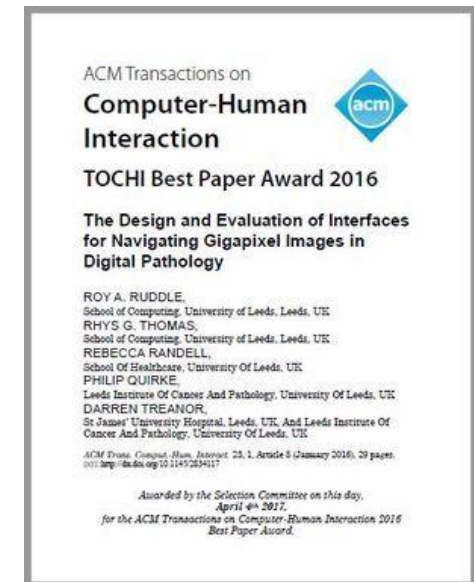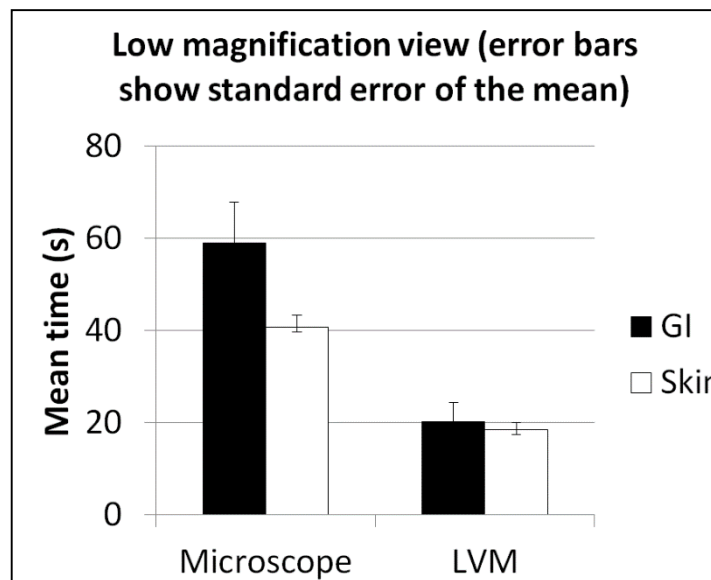- **Novel user interface**

# Evaluation: LVM vs. microscope

- **Controlled experiments**
  - **Real work (repeat diagnoses)**
  - **Participants were pathologists**
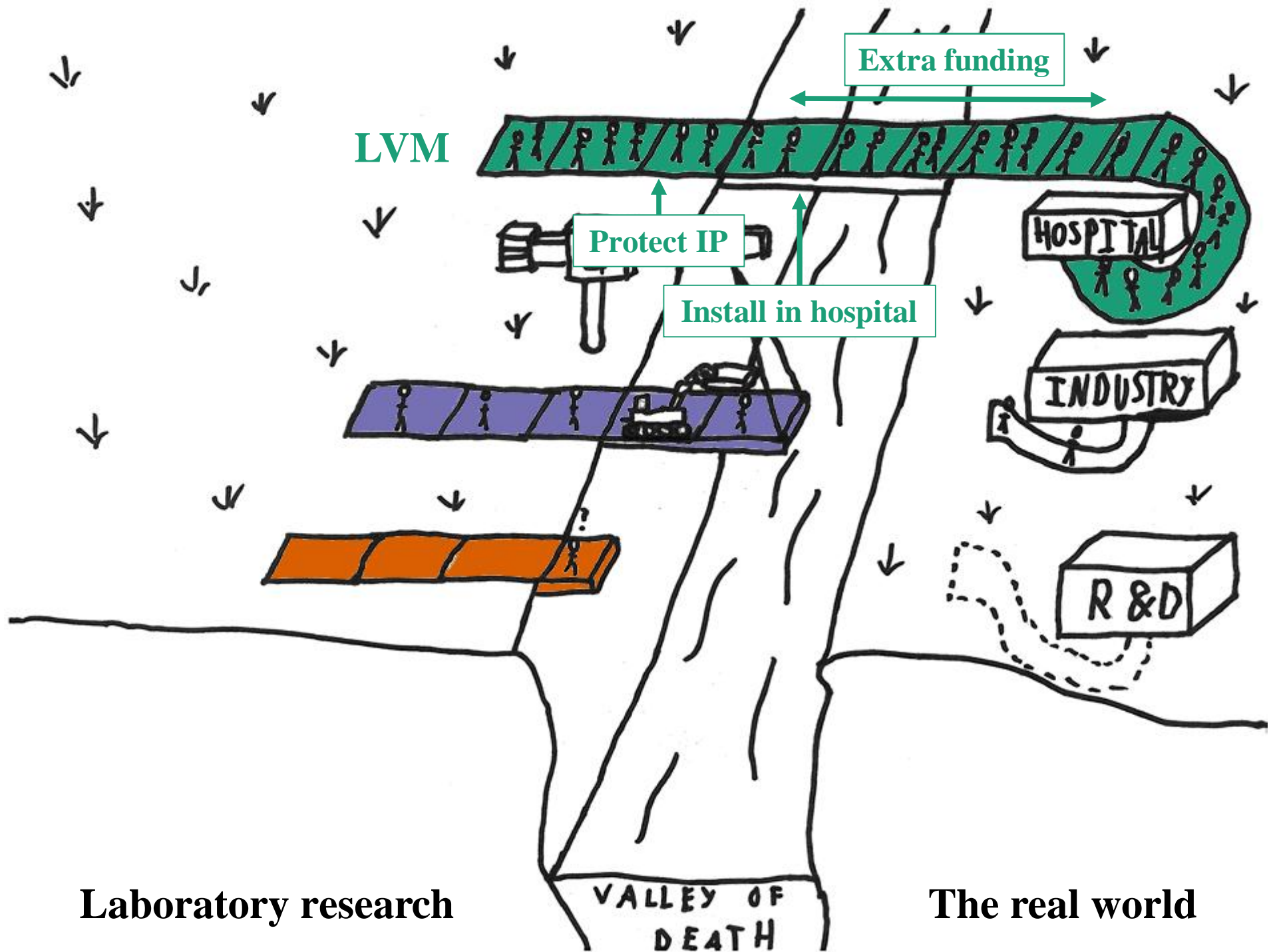    - **0.5 – 28 years experience (microscope) vs. < 1 hour (LVM)**



Diagnosis time (Error bars = 1 SD)

| **Single-slide cases** | **Long (12-25 slide) cases** | **Meta-analysis** |
|:---:|:---:|:---:|
| Randell et al. (2013). *Histopathology.* | Randell et al. (2014). *Human Pathology.* | Ruddle et al. (2016). *ACM ToCHI.* |

# What are the ingredients of the perfect project?

# Conclusions & future work

- **Generic**
  - **How can you find out what's important to users?**
  - **How can you help users imagine the benefits?**
  - **What are the ingredients of the perfect project?**
- **Visualization is fantastic for revealing anomalies**
  - **Unrealised Powerwall potential (4k is a commodity)**
- **Open the black box by visualizing detail in context**
  - **Visualization for pipeline design[1]**
  - **Visualization for machine learning (Vis4ML)[2]**
- **User interfaces**
  - **Minimise the cost ("… achieved something in minutes that would previously taken days"[3])**

[1]von Landesberger et al. (2017). *IEEE TVCG*.
[2]Sacha et al. (2018). *IEEE TVCG*.
[3]Harrison et al. (2017). *IEEE TVCG*.

# Acknowledgements

| Project | Collaborators |
|---|---|
| LVM | Rhys Thomas, Rebecca Randell, Phil Quirke, Darren Treanor (LTHT) |
| Orchestral | Peter Sondergeld, Waleed Fateen, Phil Quirke Darren Treanor (LTHT) |
| QuantiCode | Georgios Aivaliotis, Mark Birkin, Justin Keen, Kevin Macnish, Alex Markham, Chris Megone, Muhammad Adnan, Richard Kavanagh, Anna Palczewska, Jan Palczewski.<br><br>aql, Bradford Institute for Health Research, Consumerdata, Leeds City Council, Leeds Informatics Board, NHS Digital, Sainsbury's |
| Artwork | Sebastian Ruddle |

**Questions?**