

## Understanding the ‘intensive’ in ‘data intensive research’:

Data flows in next generation sequencing and environmental networked sensors

Ruth McNally  
Senior Research Fellow  
ESRC Cesagen  
Lancaster University

Co-authors:

**Adrian Mackenzie**, Cesagen, Lancaster

**Jennifer Tomomitsu**, (formerly Cesagen,  
Lancaster)

**Allison Hui**, David C. Lam Institute for East-  
West Studies, Hong Kong Baptist University  
(also formerly Cesagen, Lancaster).

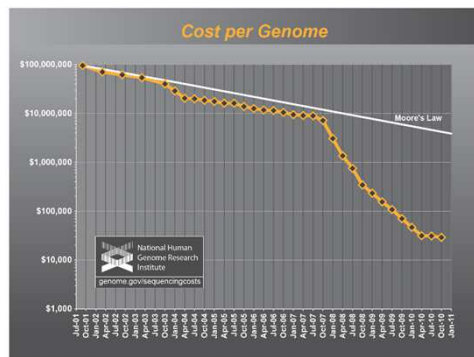
This research was undertaken with support  
from the e-Science Institute, Edinburgh

# Our motivation

## DNA Sequencing Caught in Deluge of Data

New York Times 30 Nov 2011

‘The lower cost, along with increasing speed, has led to a huge increase in how much sequencing data is being produced. World capacity is now 13 quadrillion DNA bases a year, an amount that would fill a stack of DVDs two miles high, ....’



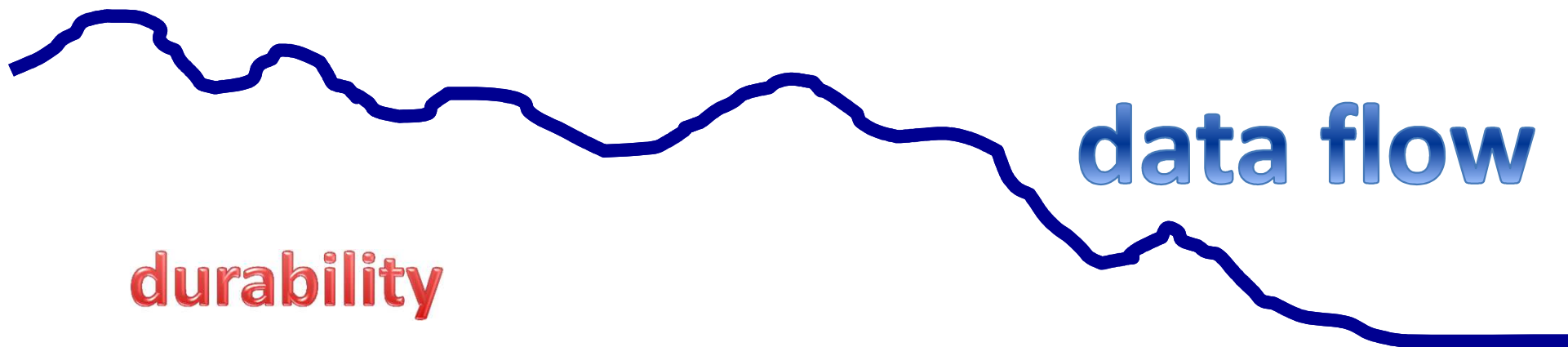
**topography**

**data flow**

**durability**

**replicability**

**metrology**





## Commoditisation of sequence - Personal sequencers

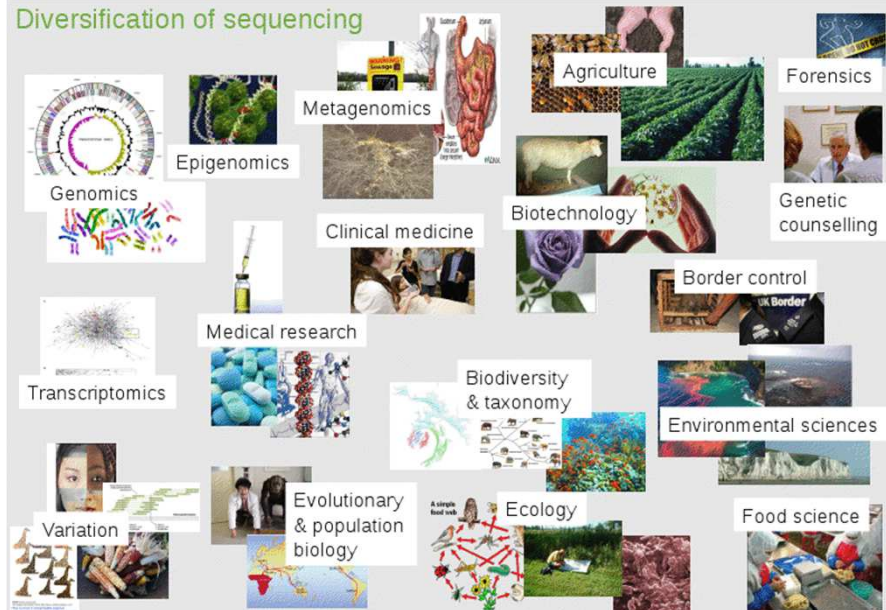


Faster benchtop machines → results in a day  
More expensive  
Diagnostics 10Mb-1Gb

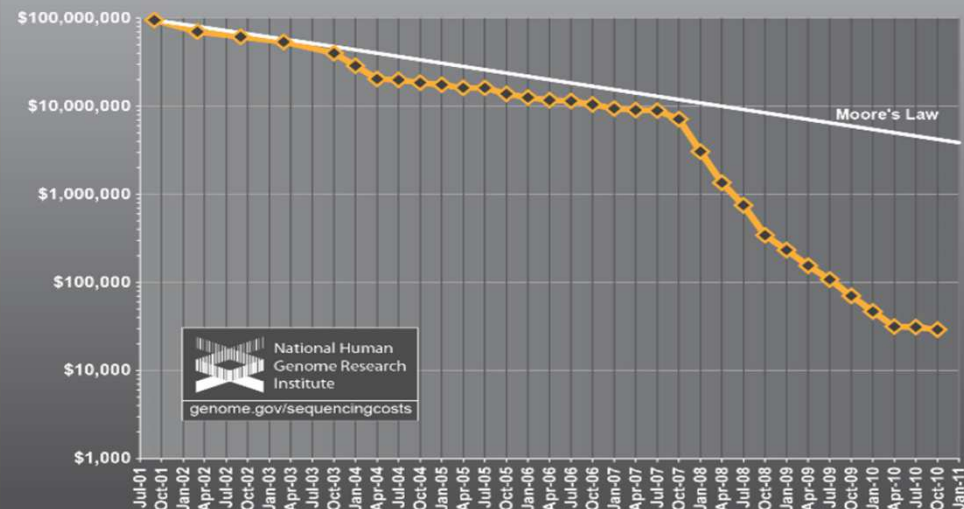
## Next Generation Sequencing

1622 total machines; inc., 712 in USA, 199 in China and 132 in UK.

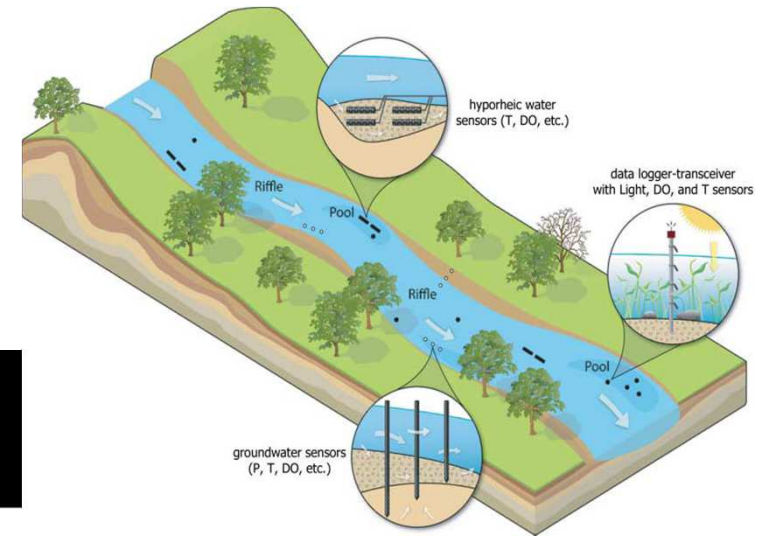
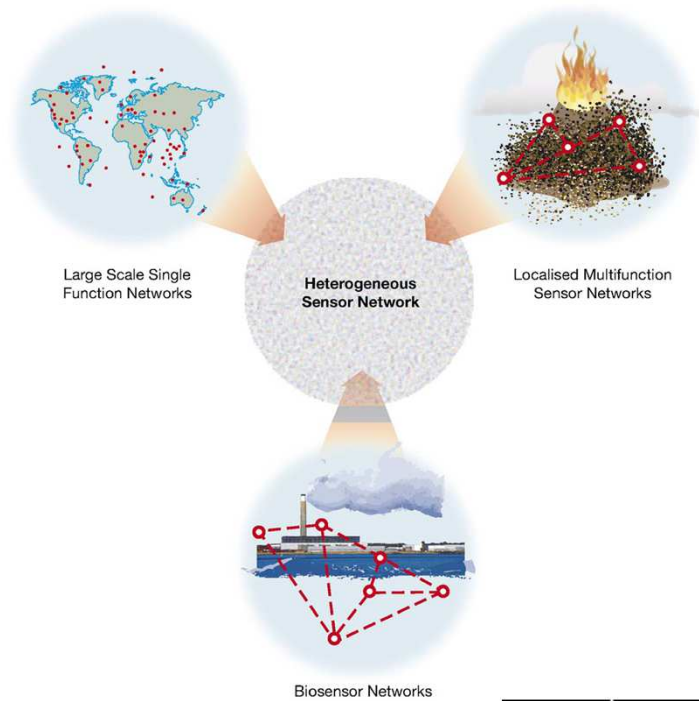
### Diversification of sequencing



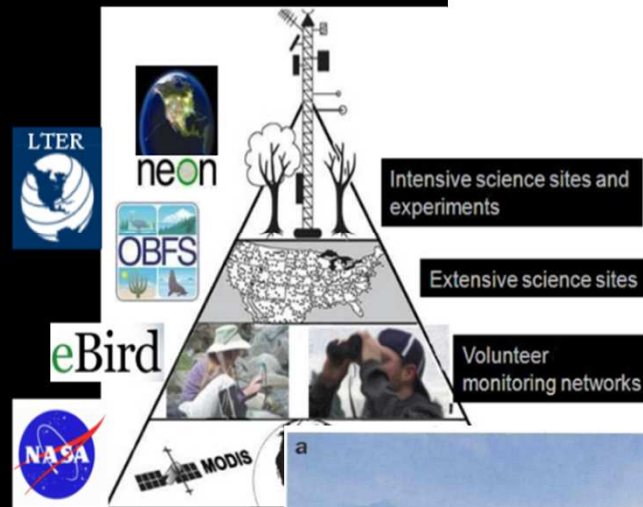
### Cost per Genome



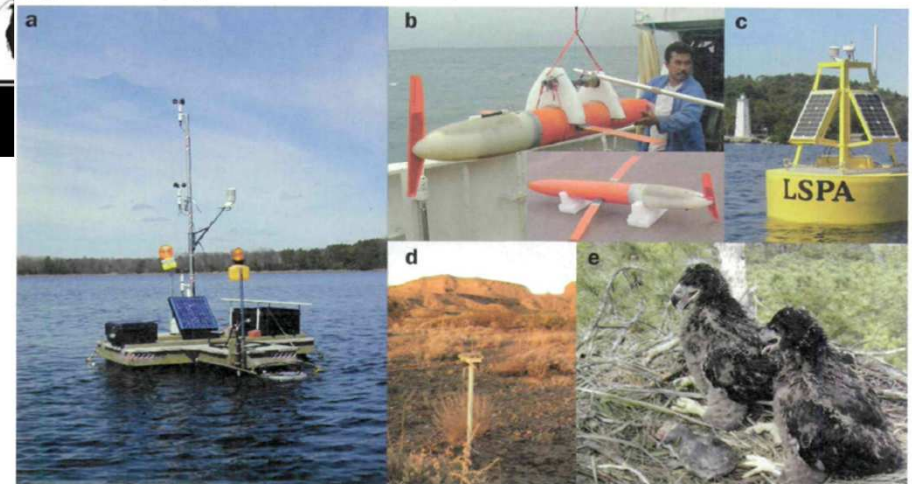
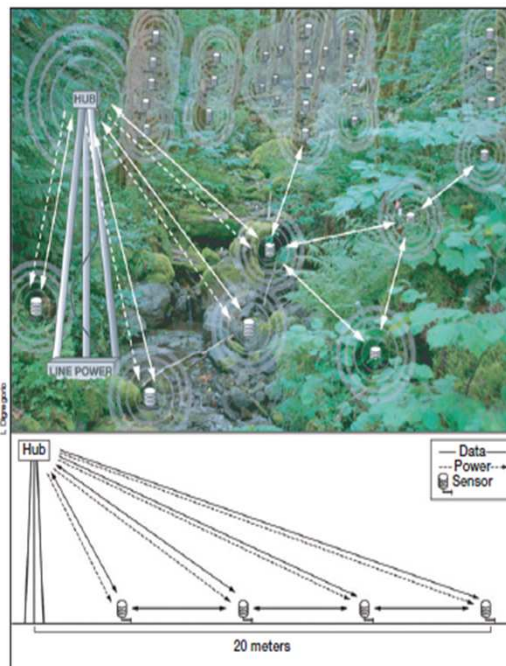




## The Earth Observation Network



## Embedded Networked Sensors

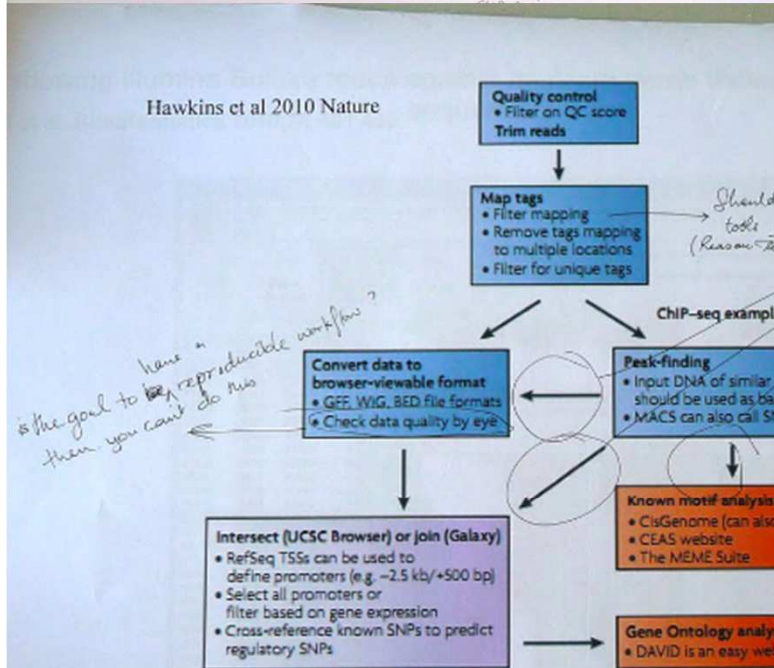
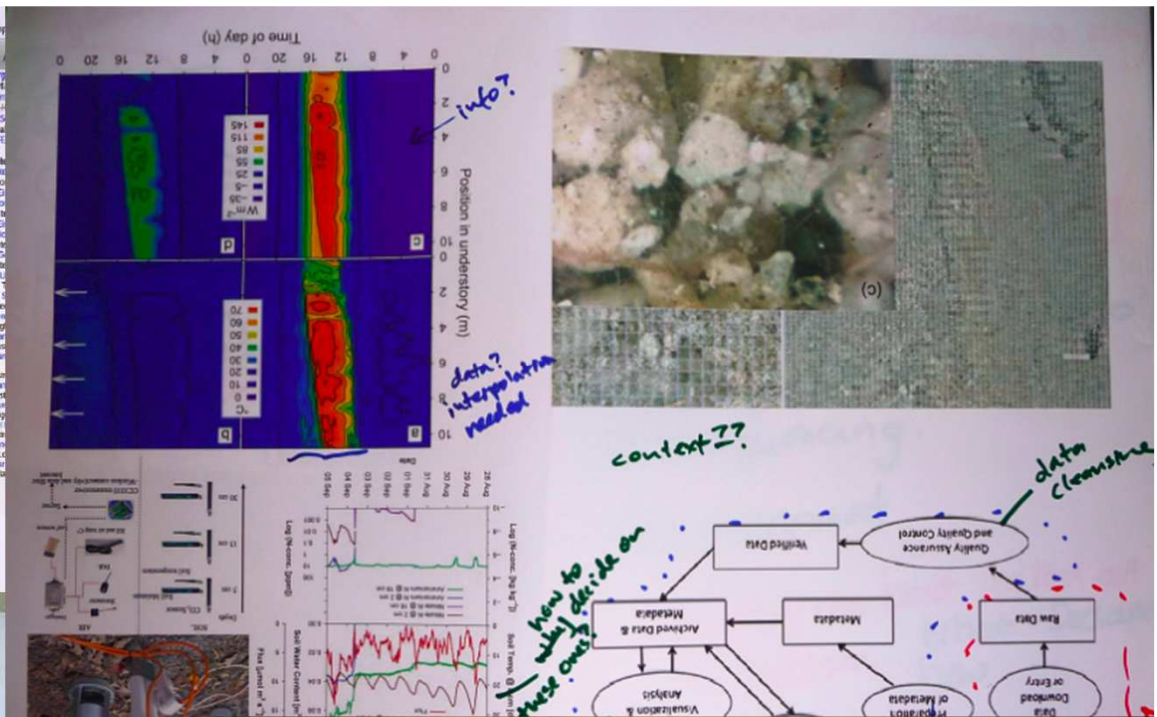




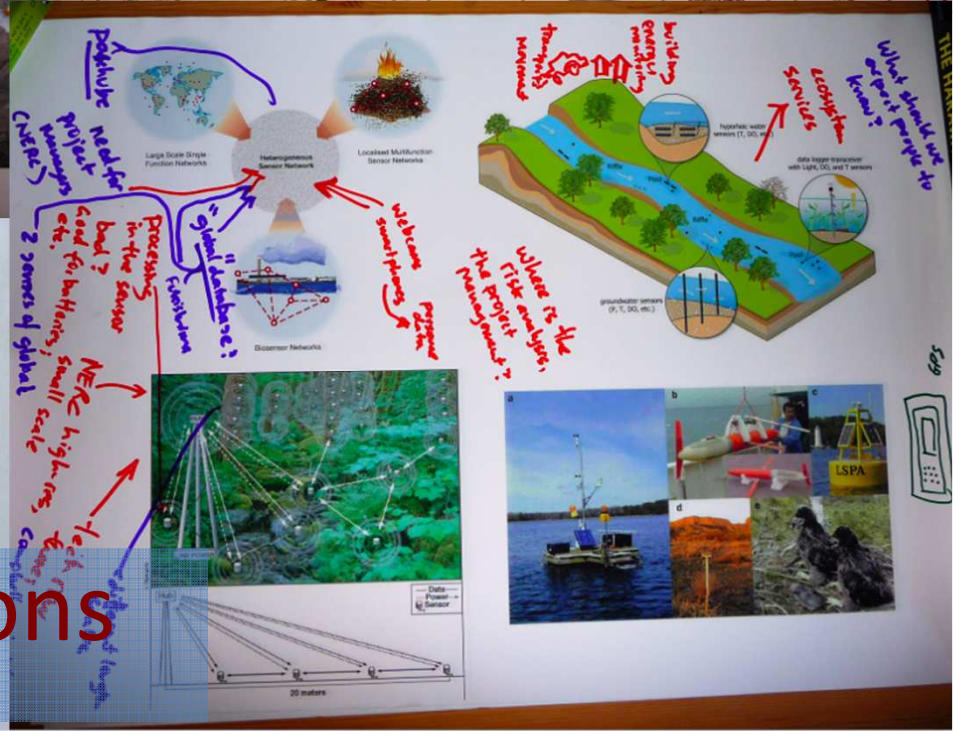


Instrumented workshops





# Collective reannotations





# Coded transcript from our own archive of digital data

282	Some of the delays in the 1000 genomes was getting a more reliable consent. Many delays in 1000 GP have to do with sample collection – the blood. You can't take 1000 GP consent back. Participants cannot change their minds.	Duration of consent – must be irrevocable
283	A big disconnect between what we see of value and what they think can get back. Disconnect between value as seen by us and as seen by their communities.	Different distributions of value
284	Some communities did say yes but then there were local diseases that they could get something back for. Mark: some did see local inheritance diseases and so said yes. Communities for help.	
285	In Africa, some of the samples came from the region in <a href="#">Kenya</a> . <a href="#">Mende</a> are so	
286	What about all the other human genomes that have been sequenced? What about the other genomes outside of 1000GP?	
287	If they are BAM, no problem. The most sequenced individual is a child whose parents who have been done.	
288	Can <a href="#">Khan</a> – not as good his because his is diploid. We are trying to get variants phasing. And getting parents with children.	
289	BLANK – 1000 genomes update – HIG data processing pipeline	
290	There has only been to put a short while!	
291	MCCARTHY 2 The Sanger perspective. Data processing and data analysis. Data production: sequencing QC Data processing: mapping, time or event, merging Data analysis: variant calling – <a href="#">SNPs</a> , <a href="#">indels</a> , other structural variants	FLOWCHART: data flows. How does this map onto Laura's Slide 12?
292	BAM production for what is released to DCC then to rest of group and then the world at large	What is DCC? – Data collection centre
293	MCCARTHY 3: METADATA TRACKING <a href="#">Metadata</a> tracking The DCC releases a monthly sequence index file. Update an internal <a href="#">genotyping</a> tracking database. Keep in sync with internal storage hierarchy. Lanes -> Library -> Platform -> Sample -> Population -> Project API to access the <a href="#">sequencing</a> database DCC is looking to use a sequence index file - we update our own <a href="#">genotyping</a> tracking database, to sync with sequence index file. Go through update <a href="#">bam</a> files for mistakes. Internal pipelines have access to the tracking database Shows an old database schema diagram	
294	MCCARTHY 4: LANE LEVEL OPERATIONS Having got <a href="#">bam</a> files, we put them into our mapping pipeline that is parallel. We have multiple methods for calling different alignments. We have other mappings ready to go. We merge to get <a href="#">bam</a> file and record <a href="#">bam</a> on the way. Then we go into <a href="#">SAMtools</a> <a href="#">sort</a> - give it a list of known <a href="#">bam</a> files, and it merges, sort and re-aligns. All of this is done on each lane. Monthly releases keep it going. After the mapping they do a merge and get in <a href="#">bam</a> the out. Merge in -> -> MAPPING -> <a href="#">bam</a> out <a href="#">bam</a> P -> SAM IMPROVEMENT -> <a href="#">bam</a> out Process and then start building a release.	THIS MAPS ONTO LAURA'S SLIDE 14: QUALITY CONTROL
295	MCCARTHY 5: BUILDING A RELEASE At a certain point, there will be a freeze - 23 November. So we freeze and start building a release. Multiple lanes and maybe even multiple libraries for an individual. Merge all lanes into a common library. MCCARTHY 6: M7, M8 - Building a release. We start with each lane. We then merge all lanes in a common library. A common platform BAM. M7, M8, M9: So for each individual, we are going to produce data for release to the DCC. M11, M12, M13, M14: On a per chromosome basis, we split <a href="#">bam</a> for all mapped reads, and a <a href="#">bam</a> for all unmapped <a href="#">bam</a> , and initially Chromosome 20 <a href="#">bam</a> because it is so small, and gives them something to do.	
296	A BAM is an individual on a common platform. If a different platform for the same individual is used, then a different BAM. <a href="#">bam</a> <a href="#">bam</a> into chromosomes and a <a href="#">bam</a> for all the mapped reads and a <a href="#">bam</a> for the bits that didn't map to the ref genome. Then a <a href="#">bam</a> release to the DCC.	
297	M15: In Phase 1 release, 1094 individuals (all 2x; 1000 > 3x). Of the 25 Tn in 2500 FASTQ FILES, about 1 RTN was generated at Sanger. BAM release data. Phase 1 and phase 1 finished. Phase 1 is 1094 individuals. Sanger and TGM in Texas. BAMS also coming from TGM and from <a href="#">genotyping</a> and solid VCF format.	metrics: data on: number of individuals, coverage, number of lanes, number of data diversity on image file <a href="#">bam</a> slide 7 - metrics on data flow in image
298	M16: With all the <a href="#">bam</a> in, we are going to start looking for variants VARIANT CALLING: VCF FORMAT: VARIATION CALL FORMAT Fully adopted by 1000G group as interchange format for variant calls <a href="#">SNPs</a> , <a href="#">indels</a> , and recently <a href="#">SVs</a> . <a href="#">Genotyping</a> calls for all samples Association of variants via web-defined to go VCF <a href="#">API</a> and tools via <a href="#">http://vcftools.sourceforge.net</a> Scaling issues with VCF - BCF format in development	IMAGE OF VCF FORMAT - HAS BECOME A STANDARD IN IUG GROUP Data diversity - VCF standard SEE LAURA'S SLIDE 13
	M17-M20: VARIANT CALLING: CHALLENGES Variant calling: challenges Traditionally BAM files have been produced per sample with all of the lanes/libraries merged Lanes -> Library -> Platform -> Sample (1 per individual) Problem: population based SNP calling needs to be aware of the reads Issues in variant calling across time and across the ref	

# Durability 1: The timing of data flow has to be synchronised with domain specific temporal dynamics

‘time can be wasted’ taking advice on experimental redesign that delays the start of data flow.  
(NGS)

‘In CENS [a project], big issue was time stamps – notoriously bad. Sad stories about non-synched data Data sets not synched properly’ (ENS)



## Durability 2: Technical and domain scientists inhabit different 'time zones'

'There was only about 2 years when data was collected that was of use to the application scientists. The initial period all about battery life, sensors, networks. They realized in the middle that it was important to keep the human in the loop –that coincided with about 2 years of useful data for application scientists. At the end of that, the technology was mature enough the application scientists could take with them and use it. The technology people got bored at this point and moved on to doing mobile applications – they kicked environmental scientists out of the loop' (ENS)

# Durability 3: Projects change

‘Projects can change from being one type of project into another ... People who got grants to do exome capture are now going to complete genomics to get analysis’ (NGS)

‘It is the Achilles heel of every semantic integration technology that it is not robust with changes. They use the most robust one (in practice). At the moment, in terms of reliable technology, it is not that scalable. The problem is mainly that modifications cause you to have a propagation effect on the mappings’ (ENS)



# Replicability 1: Too much - too little?

## Experimental replicability

'Short read sequencing is so cheap, it's a disposable item. It's cheaper make and analyse your own data than to download someone else's' (NGS)

.

# Replicability 2: Too much - too little?

## Code and practice

‘Bioinformaticians are doing the same things over and over again. Everyone has to continue reinventing the wheel. Rinse and repeat all over the world’ (NGS)

‘Most of these things [workflows] are moving targets – in our experience for mapping and assembly, how often do we change a version of it? Hourly seems to be the response’ (NGS)

‘I don't think we will ever get to fixed workflows. You will never get around to having to write new code for projects. The driver of that is the science. Science has to be novel and therefore cannot reuse whole system. That novelty is what makes you have to write new bits of code’ (NGS)



# Replicability 3: Scaling up - is repetition enough?

## Collaboration and enrolment

Can't do this on your own – have to have a massive team – computer scientists, engineers, domain scientists, people to keep spirits up.” (ENS)

‘You have to demonstrate it works as well as previous methods or better, and then wait for acceptance from the discipline before you go too far’ (ENS).

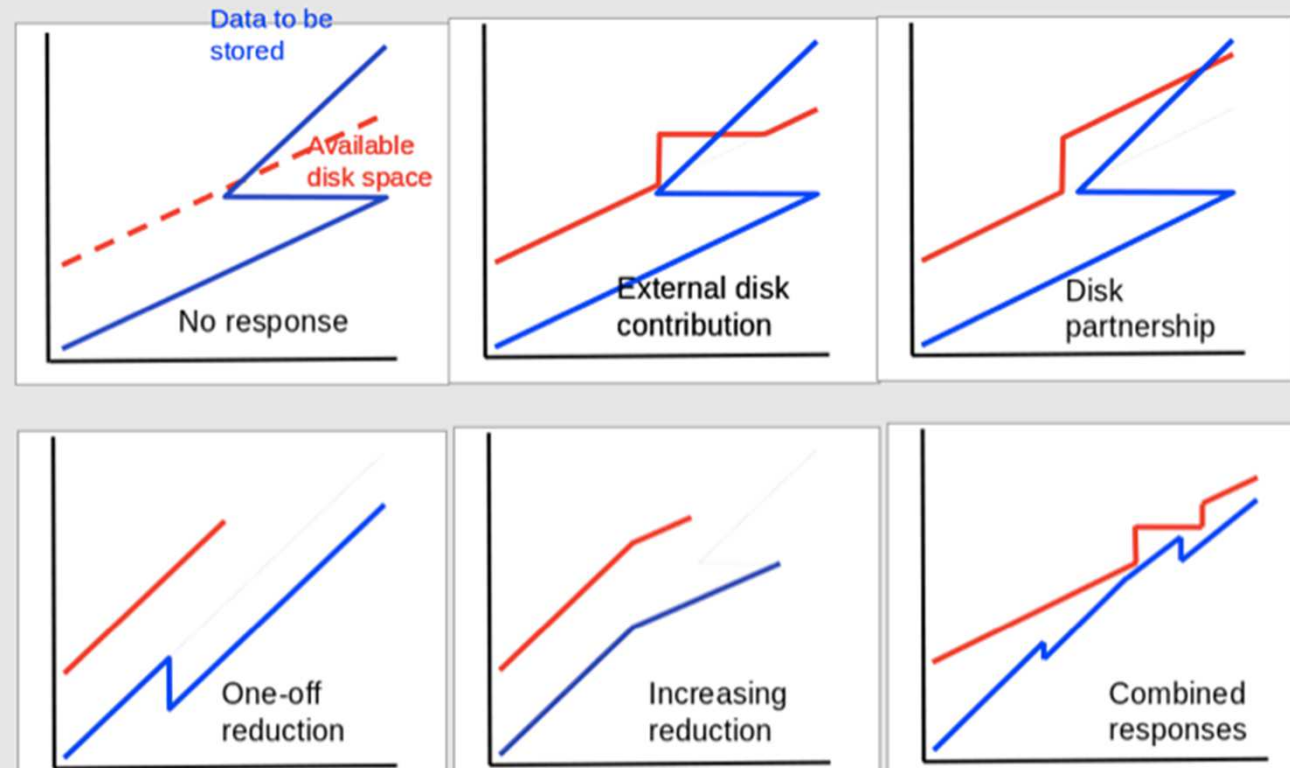
## System change

‘One of projects – eBird – global project – concept is to get volunteers to go out and using fairly standard protocols, but standard, collect their observations of birds. ... When first started, couldn't get anybody to do that. So we changed how we thought about citizen science data. Changed in 2005. Launch of eBird 2.0. Last Tuesday they collected more data than they did in 2004’ (ENS)

# Metrology 1

Keeping  
within the  
curve

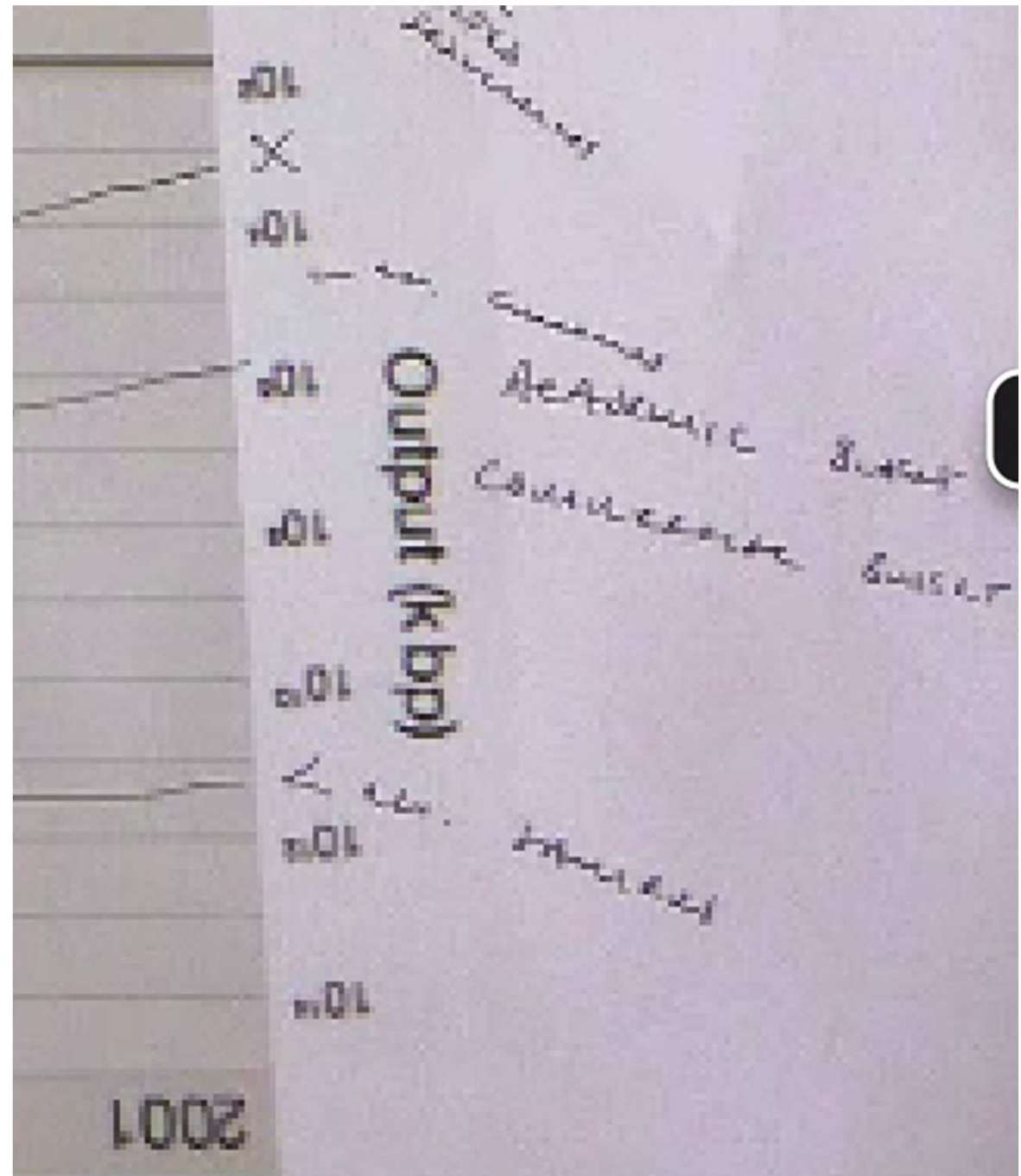
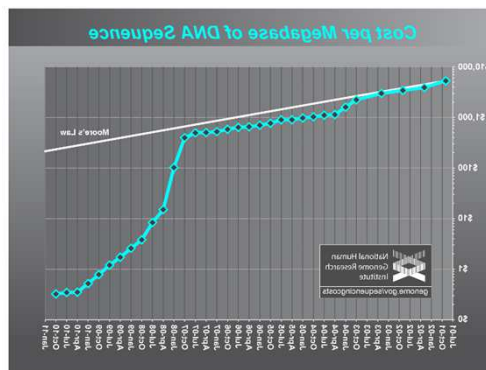
## Strategy for handling data growth



‘If we have a fixed annual spend, we will be able to grow our disk space exponentially. The data is also exponential. The risk is that if we do nothing, two lines cross. I'm going to show an exercise in keeping the blue line below the red line. Either partner on disk space. The kind of compression that will deflect this curve. There are features that will allow that to happen. The strategy is to have a set of responses to hand that we can deploy when we need to. Apply strategies of data reduction, judicious and community informed’ (NGS)

# Metrology 2

missing metrics



# Metrology 3: Novel metrics and data bibliometrics

‘Recently an ecologist determined you could more accurately determine the onset of spring through public webcams using green divided by blue than by using remote sensing data’ (ENS)

‘Is there any benefit to having standards? You get cited more if you cite ArrayExpress. Look at ProteoRED MIAPE satisfaction survey. 95% of people like MIAPE. Papers with data in ArrayExpress get cited more than equivalent papers that don’t have data in ArrayExpress’ (NGS)



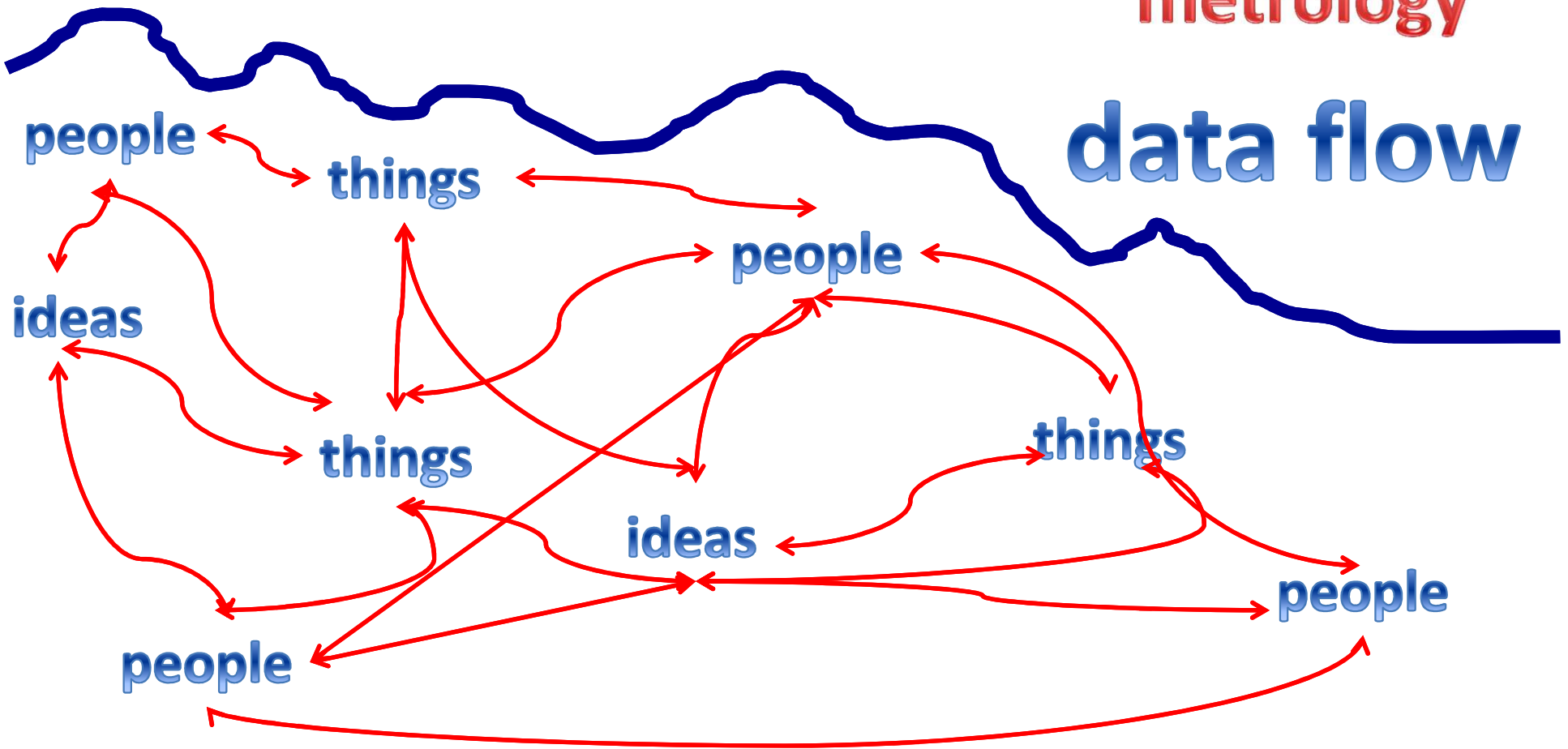
durability

replicability

metrology

topography

data flow



THANK YOU FOR YOUR ATTENTION

This research was undertaken with support from the e-Science Institute, Edinburgh (see <http://www.esi.ac.uk/research-themes/20>). The support of the Economic and Social Research Council (ESRC) is gratefully acknowledged. The work presented forms part of the programme of the ESRC Genomics Network at Cesagen.

