# CONFESSIONS OF A DATA WRANGLER

Dr Carina Kemp

Director of eResearch

Australian Academic and Research Network

IS A JOURNEY (OF DISCOVERY)
NOT A DESTINATION

What is a data wrangler?

Why are we forced to be wrangers?

Technology?
Friend or enemy…

WHAT does this all mean?

# WHAT IS A DATA WRANGLING

**https://en.wikipedia.org/wiki/Data_wrangling**

**Data wrangling**, sometimes referred to as **data munging**, is the process of transforming and mapping data from one "raw" data form into another format with the intent of making it more appropriate and valuable for a variety of downstream purposes such as analytics. A **data wrangler** is a person who performs these transformation operations.

This may include further munging, data visualization, data aggregation, training a statistical model, as well as many other potential uses. Data munging as a process typically follows a set of general steps which begin with extracting the data in a raw form from the data source, "munging" the raw data using algorithms (e.g. sorting) or parsing the data into predefined data structures, and finally depositing the resulting content into a data sink for storage and future use.

# ALL SCIENTISTS ARE DATA WRANGLER...

# CONFESSIONS OF A DATA HOARDER?

# WHY DO I WRANGLE?



data-swamp
interoperability
collaborate
usable
flexibility
lazy
workflow
software
proprietary
open-source
longtail
reusable
Wrangling?
formats
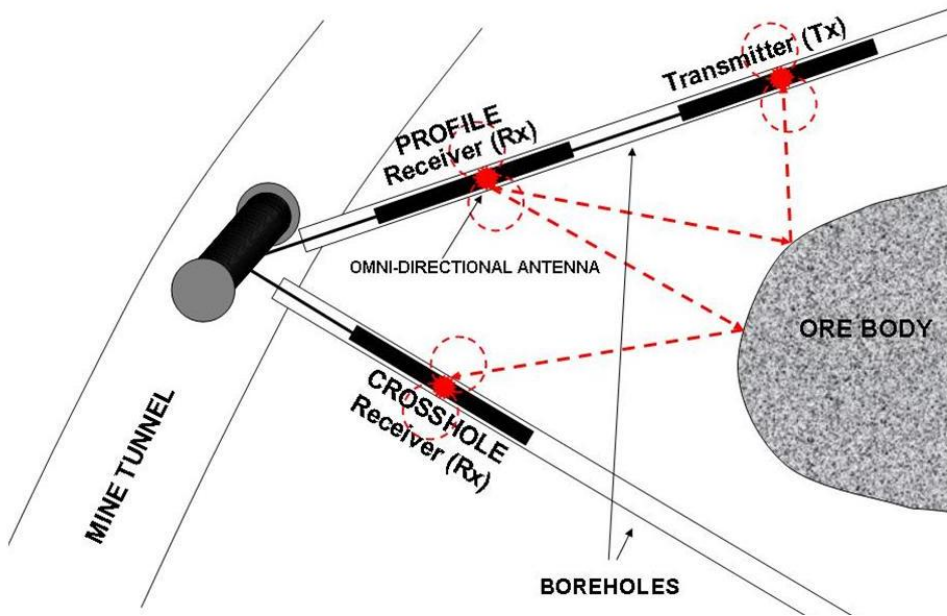exchange

# WHAT DO I WRANGLE?

# THE DATA WRANGLING BEGAN…

aarnet

# "MINE-SCALE THREE DIMENSIONAL BOREHOLE RADAR (BHR) IMAGING"

What is Borehole Radar?

The borehole radar system can be deployed by winch or on the drill rods similar to a gyro survey
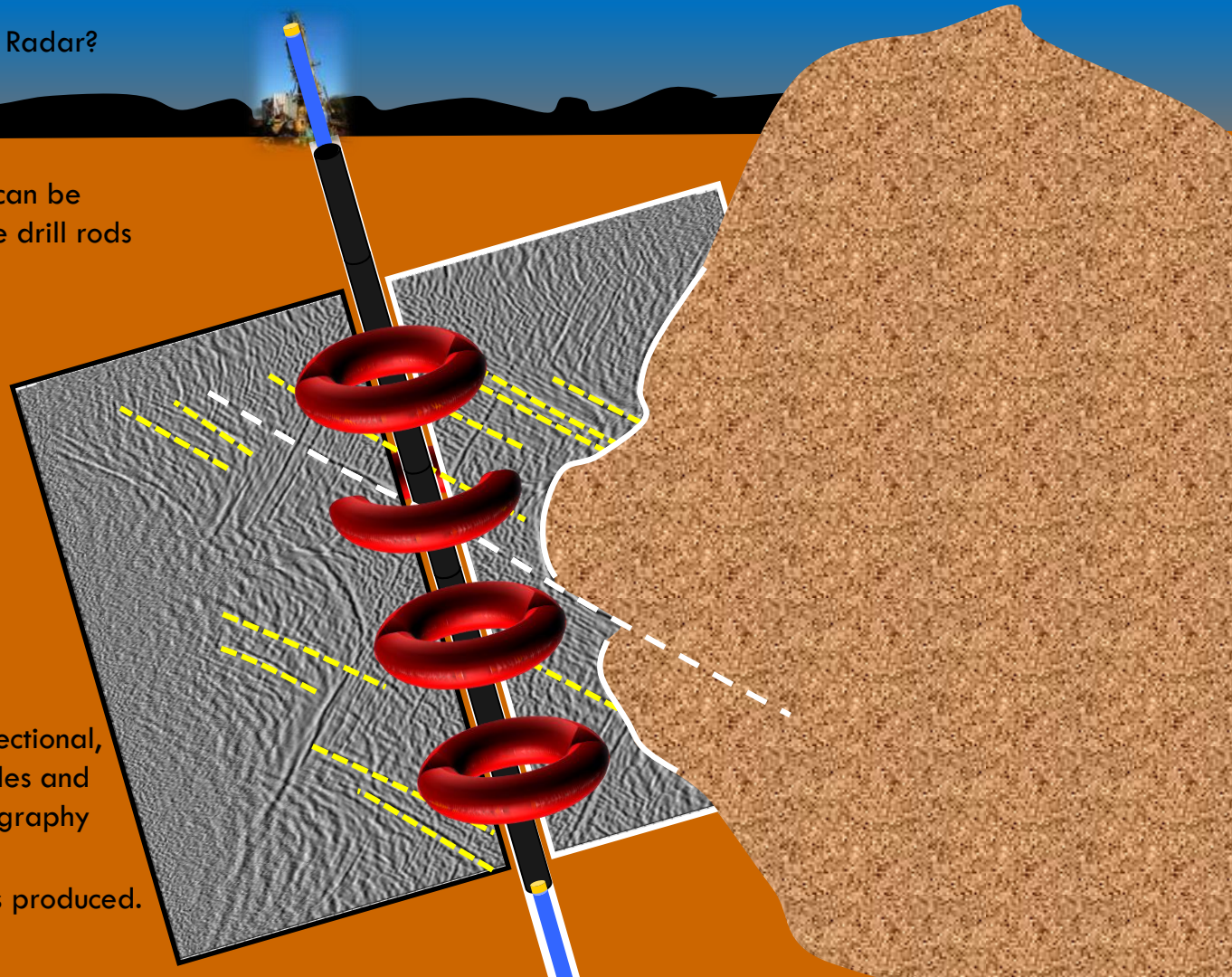
Data is acquired continuously as the rods are pulled and the radar ascends the drillhole **Signal is sent radially outwards into the surrounding rock** The radar images the rock surrounding the drillhole.
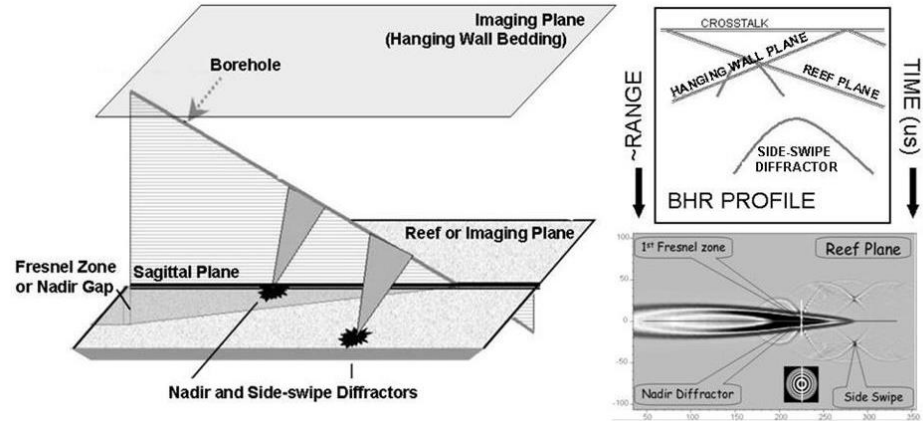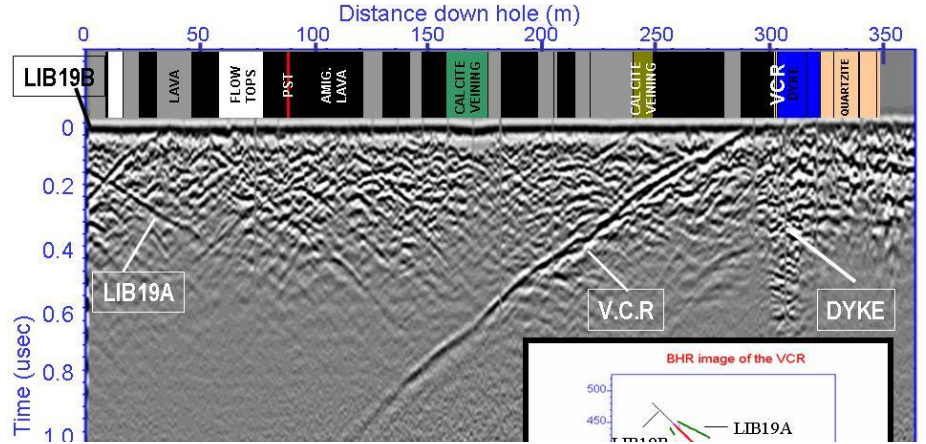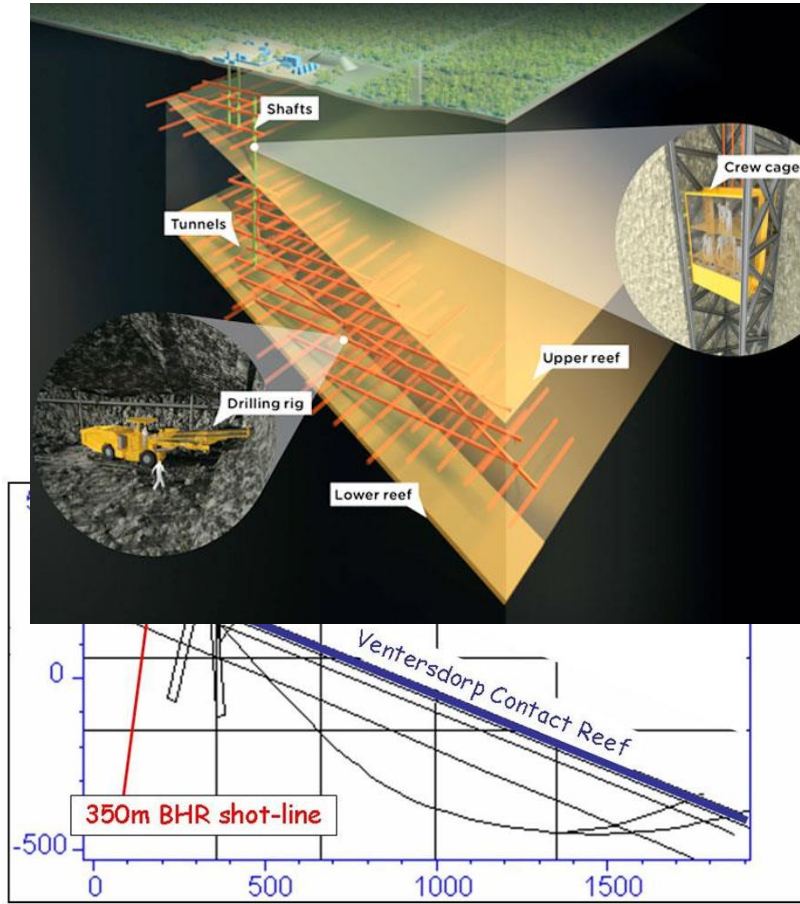
The radar is not directional, Neighboring drillholes and knowledge of stratigraphy aids interpretation.
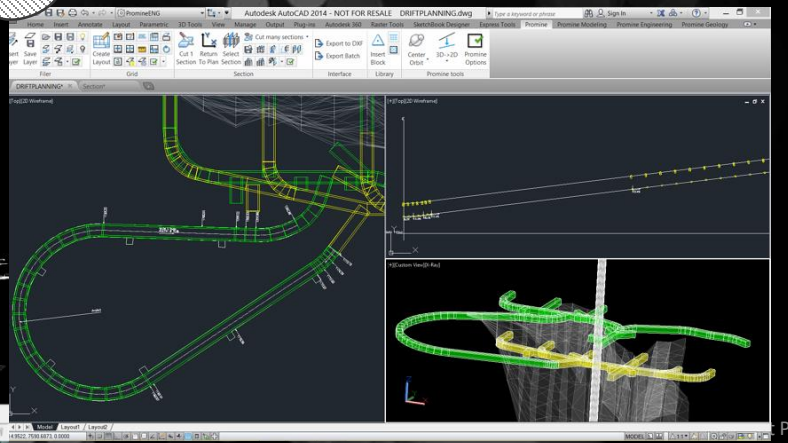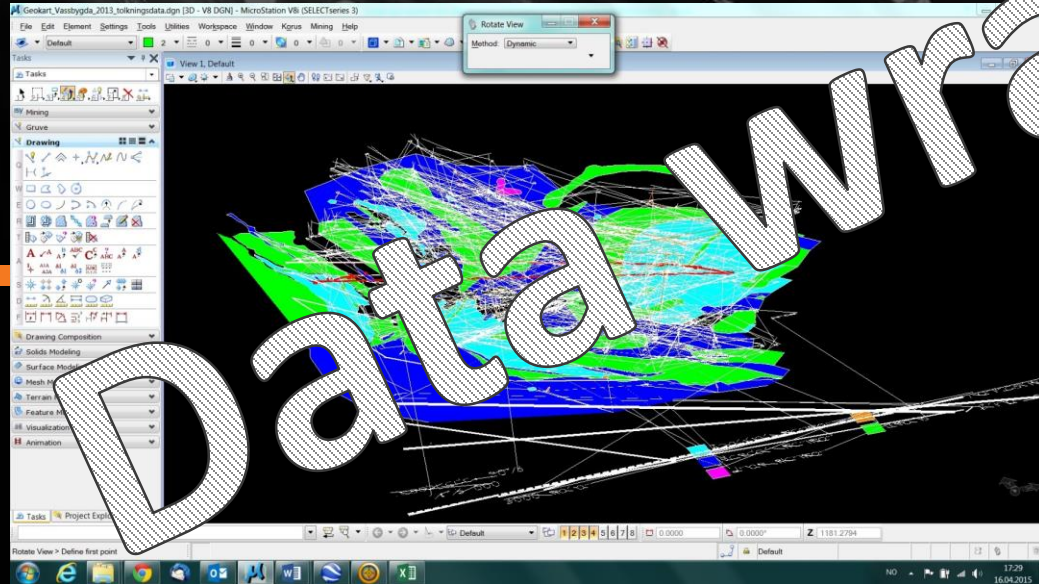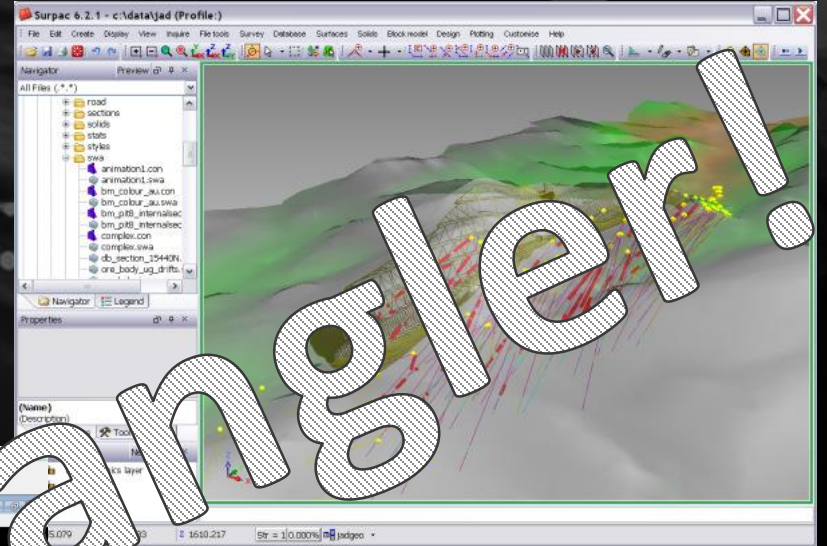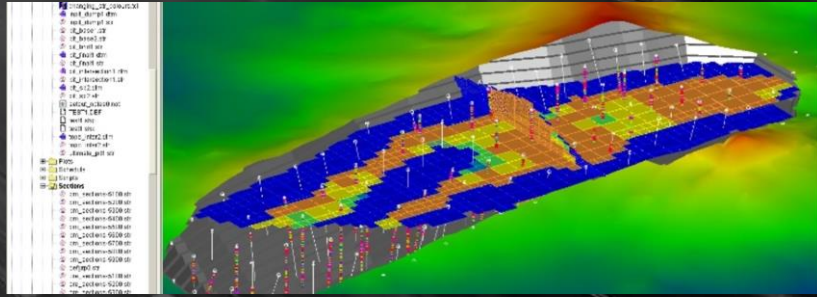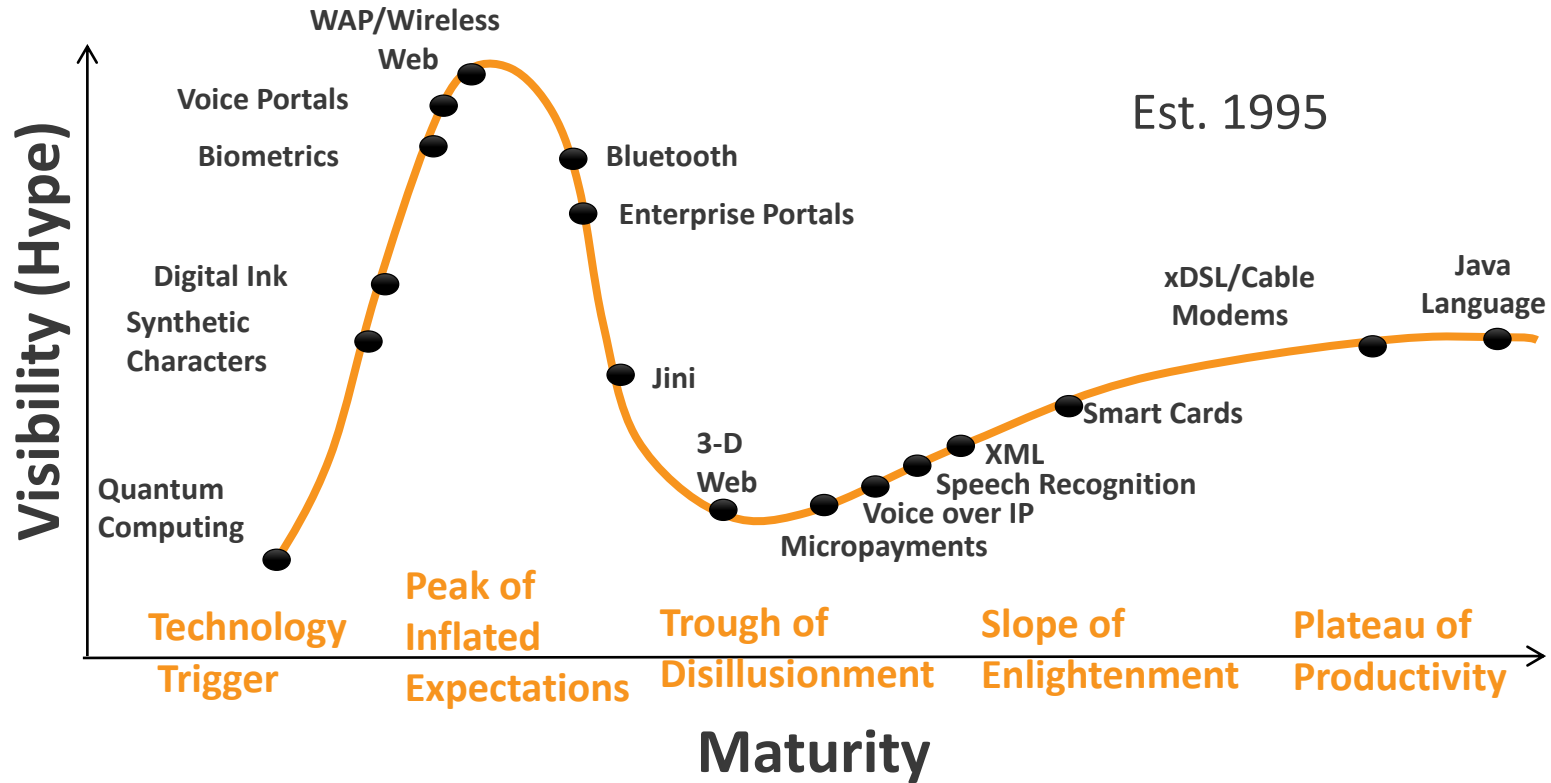
Final interpretation is produced.

# THE DATA...

THE SOFTWARE...

Data wrangler!

# TECH JOURNEY: FRIEND OR ENEMY…

# THE GARTNER HYPE CYCLE

# GARTNER HYPE CYCLE 2000



Gartner Hype Cycle chart — Visibility (Hype) vs Maturity

- WAP/Wireless Web
- Voice Portals
- Biometrics
- ASPs
- Webtops
- Bluetooth
- Enterprise Portals
- Digital Ink
- Synthetic Characters
- Audio Mining
- Quantum Computing
- Jini
- 3-D Web
- Micropayments
- Voice over IP
- XML
- Speech Recognition
- Smart Cards
- xDSL/Cable Modems
- Java Language

Maturity stages: Technology Trigger · Peak of Inflated Expectations · Trough of Disillusionment · Slope of Enlightenment · Plateau of Productivity



UNDERGROUND MINES
The 10 deepest

COMPILED BY RICHARD JANSEN VAN VUUREN

10 Creighton nickel mine 2.5 km
9 Great Noligwa gold mine 2.6 km
8 Kidd Creek copper & zinc mine 2.92 km
7 South Deep gold mine 2.99 km
6 Moab Khotsong gold mine 3.05 km
5 Kusasalethu gold mine 3.2 km
4 Driefontein gold mine 3.4 km
3 Savuka gold mine 3.7 km
2 TauTona gold mine 3.9 km
1 Mponeng gold mine +4 km

# GARTNER HYPE CYCLE 2001



**Tech**

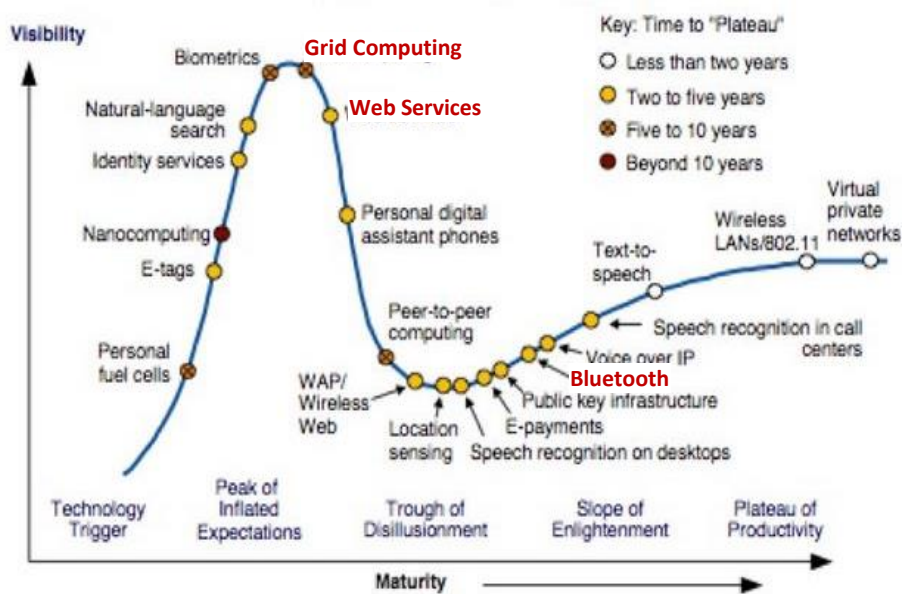Semantic Web / Web Services / WAP / Bluetooth

**Me**

Started my PhD – data collection

Some coding – Fortran / Matlab / C

**BHR /Geophysics**

Begin Testing wireless technologies

Data Standard...?

# GARTNER HYPE CYCLE 2002



**Tech**

Grid Computing
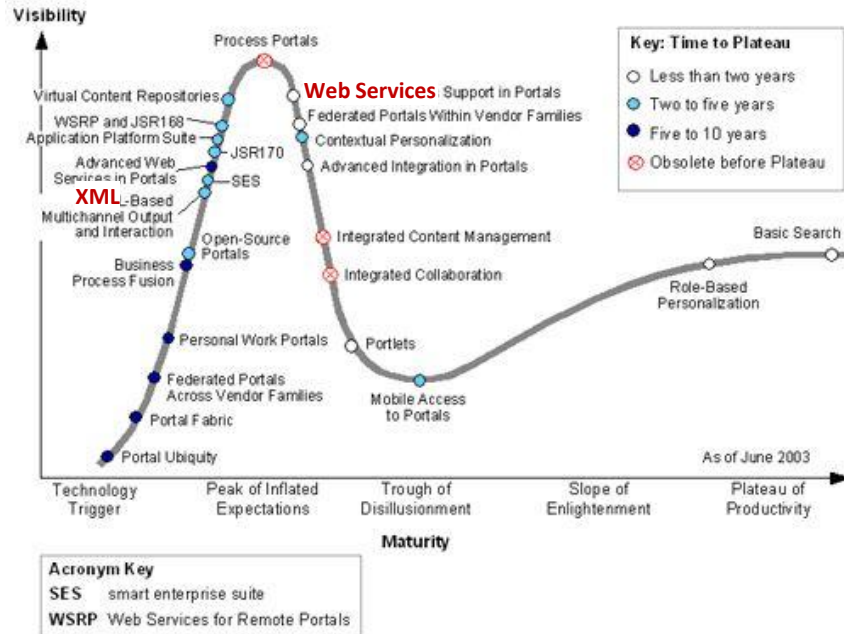
Web Services maturing

Bluetooth stabilising

**Me**

PhD Year 2 – more data collection

coding – Matlab / C

**BHR / Geophysics**

Flash memory / Bluetooth tests

# GARTNER HYPE CYCLE 2003



**Tech**

XML Based Multi-channel Output and Interaction

Web Services

**Me**   PhD….
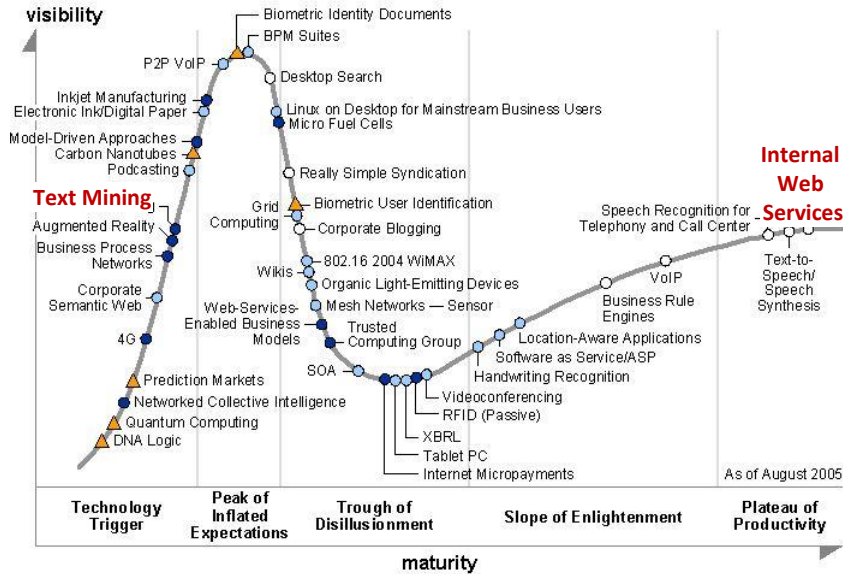
**BHR**

Single stick wireless radar tests

**Geophysics**

*ASEG-GDF2 published – A standard for point located data exchange*

Australian Society of Exploration Geophysicists

# …GARTNER HYPE CYCLE 2005



**Tech**

Text Mining…

**Me**

Started Postdoc

Data Fusion Project, early ML

**BHR / Geophysics**

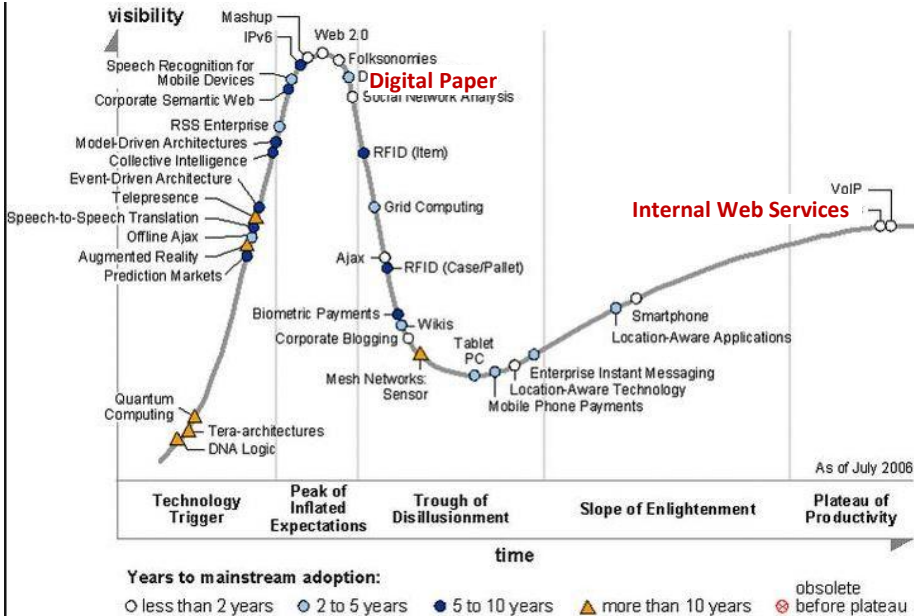Mature Bluetooth radars



Radar **+** Spacers **+** PDA **+** Drill Attachment

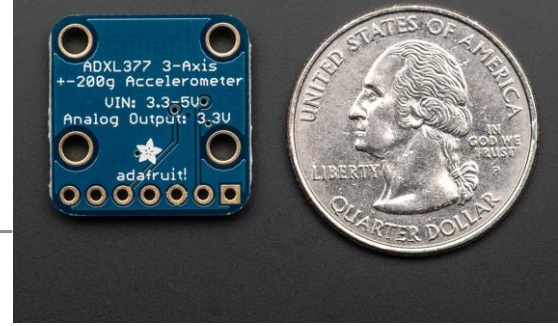# GARTNER HYPE CYCLE 2006





Economist.com

https://www.economist.com/asia/2017/03/09/the-end-of-a-mining-boom-leaves-australias-economy-surprisingly-intact

# ...GARTNER HYPE CYCLE 2009



Source: Gartner (July 2009)
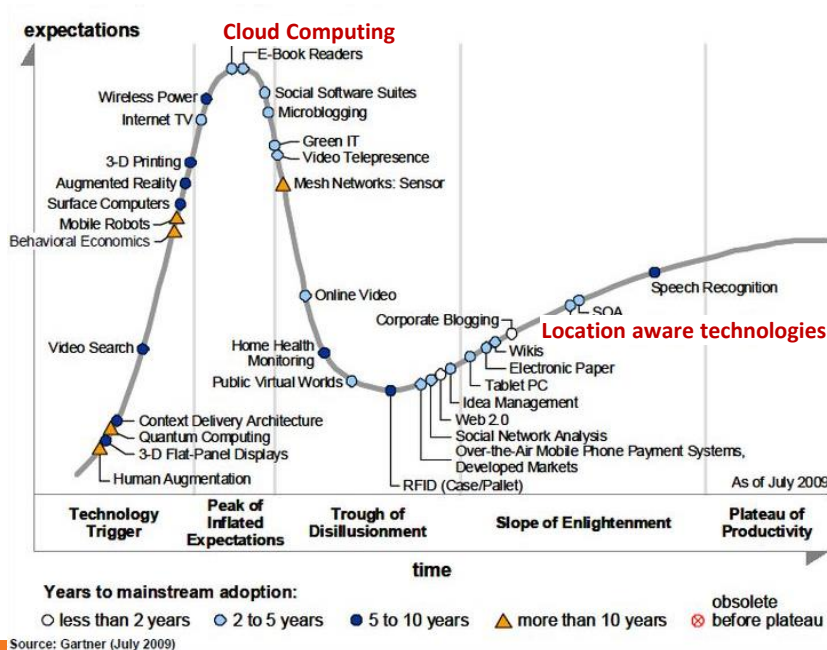
https://www.adafruit.com

**Tech**

Cloud Computing –

**Me**

BHR Business Development

GOLD / PLATINUM / NICKEL / DIAMONDS

**BHR / Geophysics**

Downhole surveying tools with accelerometers
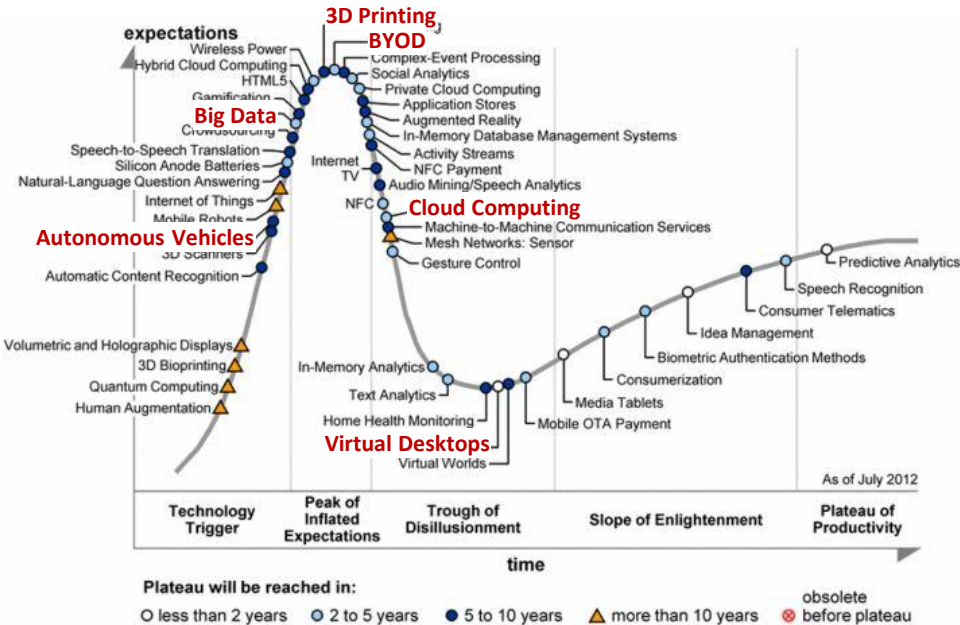
# ….GARTNER HYPE CYCLE 2012



**Tech**

Big Data is emerging

Virtual Desktops…

**Me**

Moved to Canberra

Started playing with BIGish Geophysics Data

Virtual Geophysics Exploration Laboratory (VEGL)
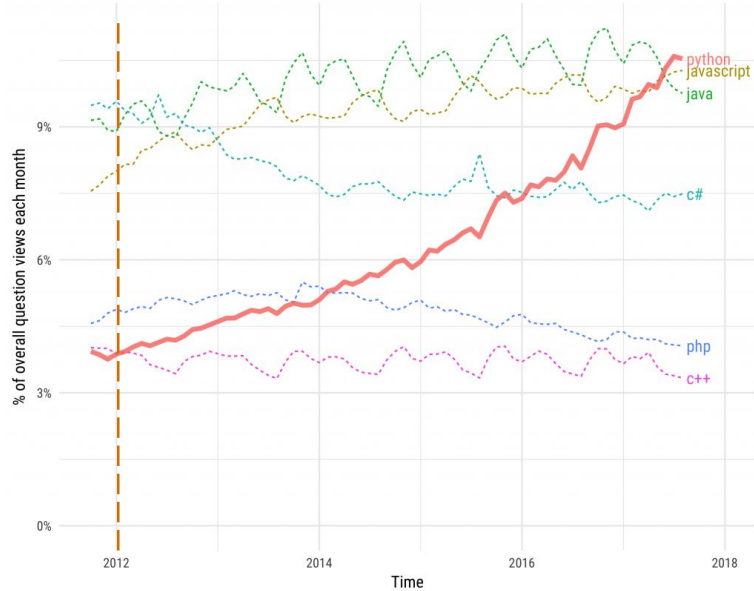
Python

~~BHR~~ / **Geophysics**

*ASEG-GDF2 still the same since 2003…*

*ASEG-ESF for Electrical Surveys published*

# ASIDE – PYTHON GROWTH…
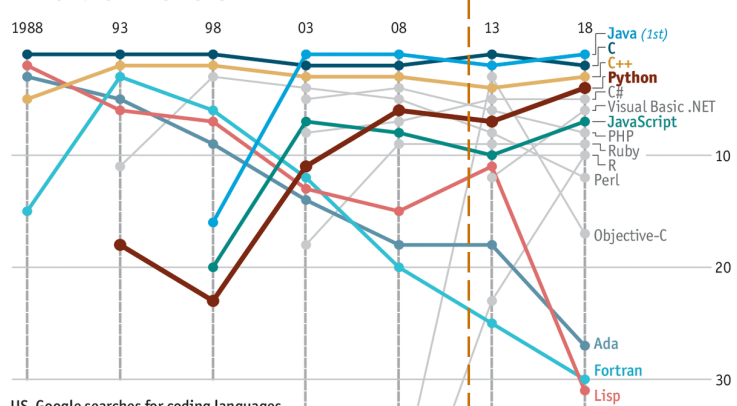


Growth of major programming languages
Based on Stack Overflow question views in World Bank high-income countries
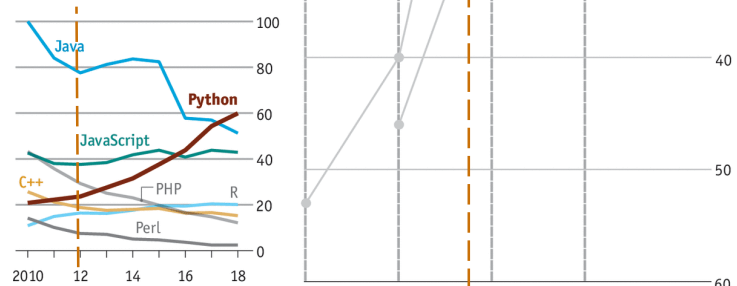


Code of conduct
Ranking of programming languages*

US, Google searches for coding languages
100=highest annual traffic for any language

Source: TIOBE, Google Trends
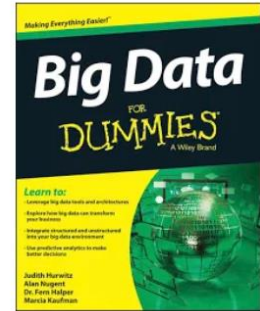
The Economist

*Ranked by global search-engine popularity

https://stackoverflow.blog/2017/09/06/incredible-growth-python/

https://www.economist.com/graphic-detail/2018/07/26/python-is-becoming-the-worlds-most-popular-coding-language

# GARTNER HYPE CYCLE 2013



**Tech**

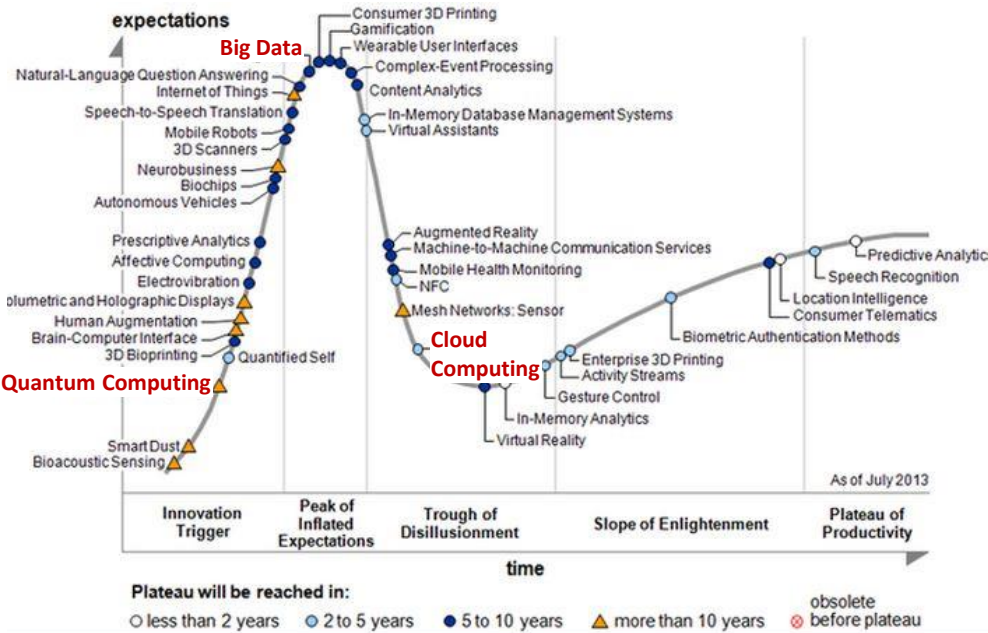Big Data hits its peak of HYPE...

Quantum Computing still emerging
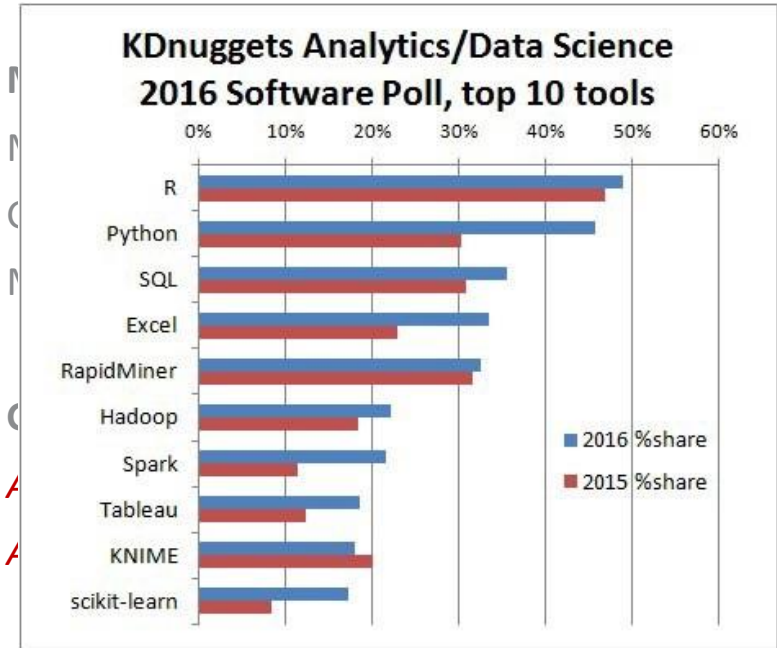
Cloud Computing in the dumps..

**Me**

Maternity Leave

Virtual Geophysics Laboratory (VGL)

# GARTNER HYPE CYCLE 2014

**Tech**

Data Science Appears

# GARTNER HYPE CYCLE 2017

**Tech**

Deep Learning Hype

Machine Learning Hype

Blockchain heading for the tr...

*Semantics and Web Services ...*

**Me**

ANVGL in the Cloud in 2015

National Geophysical Collection Maturity – NetCDF-CF

**Geophysics**

*ASEG-GDF2 still the same since 2003…*

*ASEG-ESF still the same since 2012…*

# GARTNER HYPE CYCLE 2018





Figure labels (left to right along the curve):

**Digital Twin** · **Deep Learning**

Biochips — Carbon Nanotube
Smart Workspace — IoT Platform
Brain-Computer Interface — Virtual Assistants
Autonomous Mobile Robots — Silicon Anode Batteries
Smart Robots
Deep Neural Network ASICs — **Blockchain**

**Quantum Computing**

5G — Connected Home
Self-Healing System Technology — Autonomous Driving Level 4
Conversational AI Platform
Autonomous Driving Level 5 — Mixed Reality

Edge AI
Exoskeleton

**Blockchain for Data Security**
Neuromorphic Hardware
4D Printing

Artificial General Intelligence — Smart Fabrics

Smart Dust
Flying Autonomous Vehicles — Augmented Reality
Biotech — Cultured or Artificial Tissue

Plateau will be reached in:
- less than 2 years
- 2 to 5 years
- 5 to 10 years
- more than 10 years

As of August 2018

Innovation Trigger | Peak of Inflated Expectations | Trough of Disillusionment | Slope of Enlightenment | Plateau of Productivity
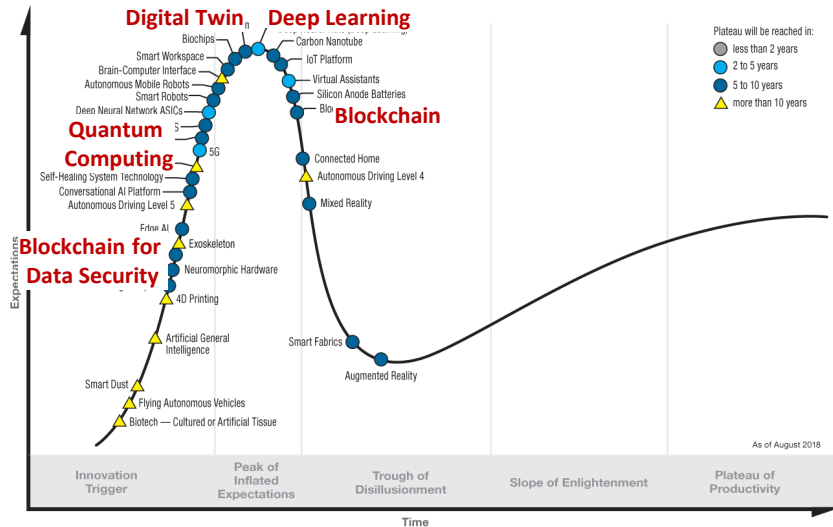
Expectations — Time

## Tech

Blockchain for Data Security

Quantum Computing is getting closer

Digital Twins??

Deep Learning

*Semantics and Web Services Matured??*

## Me

Management…  Acting CIO…

Digital Transformation

## Geophysics

*ASEG-GDF2 still the same since 2003…*

*ASEG-ESF still the same since 2012…*

WHAT CAN WE LEARN FROM
THE HYPE?

# CONFESSION 1.

**I am a bad data manager naturally**

▶

**But good practice can be learned over time and with good collaborators..**

aarnet

# CONFESSION 2.

I was a Data Wrangler and now I enable Data Wranglers..

**cloudstor**

But leveraging technologies can make it easier..

Efficient.. Repeatable..

FAIR…

# CONFESSION 3.

**Web-services to enable Machine actionable FAIR Data has been possible since 200#**

**But where are we now?**

▶

**The Geophysical Community standardized early for interoperability but not machine to machine actions…**

# CONFESSION 4.

**Change is hard** ▶ **But we must be patient and work together across our communities to promote the benefits of change.**

aarnet

# CONFESSION 5.

**Change is hard (especially in Silos)**

**Hindsight is easy**

▶

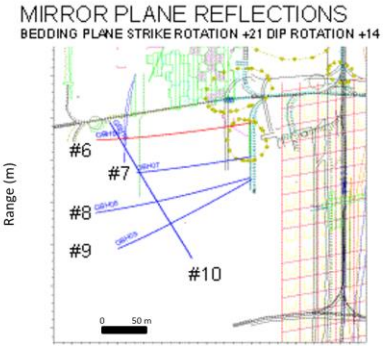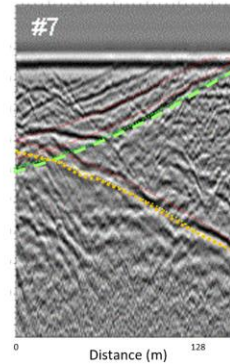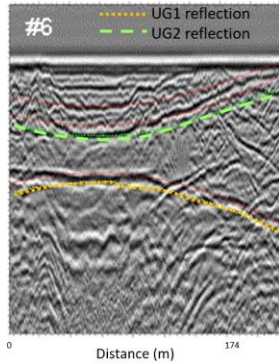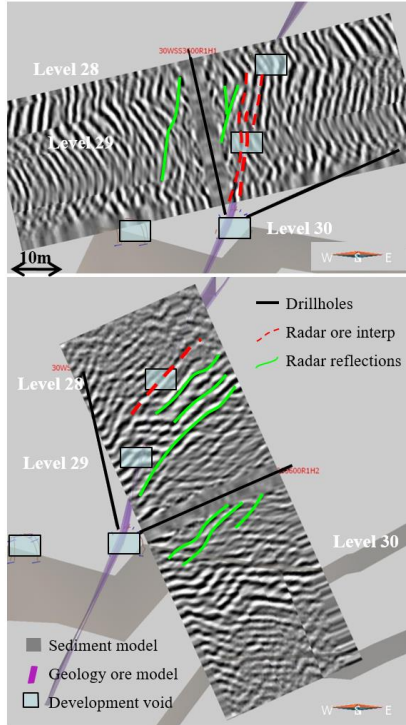**But we must be patient and work together to promote the benefits.**
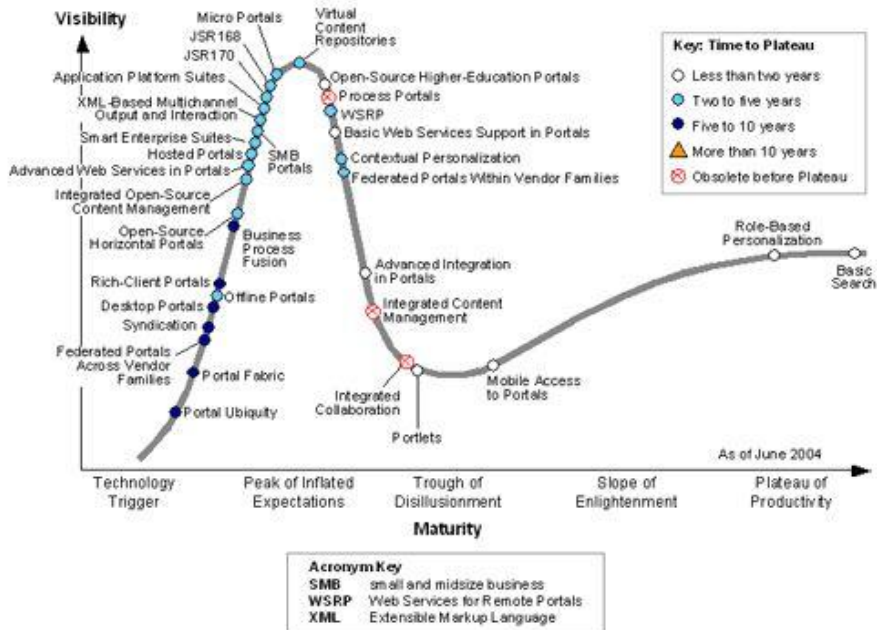
MORE RESCUED DATA…

# GOLD

# PLATINUM

THANK YOU

# GARTNER HYPE CYCLE 2004

Gartner Hype Cycle 2004

**Me**

**BHR**

**Geophysics**