# Finding Stories in Data

David Tarrant · @davetaz

# Session 2
# Finding and exploring data

# DATABLOG

## Facts are sacred

## Anyone can do it. Data journalism is the new punk

Can anyone be a data journalist? **Simon Rogers** on what we can learn from a 1977 diagram

• Another view: What data can and cannot do by Jonathan Gray

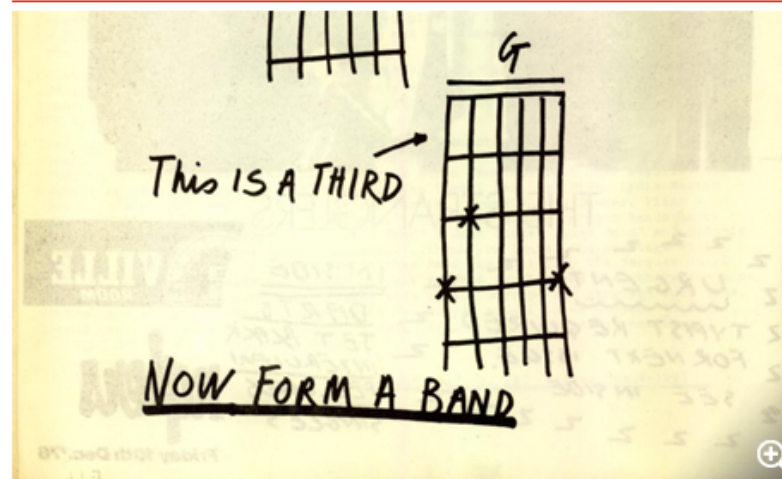Page two of Sideburns, January 1977

❝ This is a chord… this is another… this is a third. NOW FORM A BAND

Posted by
Simon Rogers
Thursday 24 May 2012
13.00 BST
theguardian.com
💬 Jump to comments (8)

Article history

**Media**
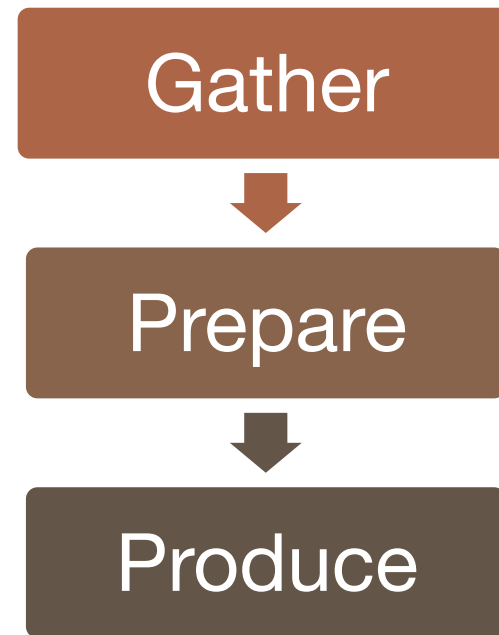Data journalism · Open journalism

**Technology**

# The data percolator

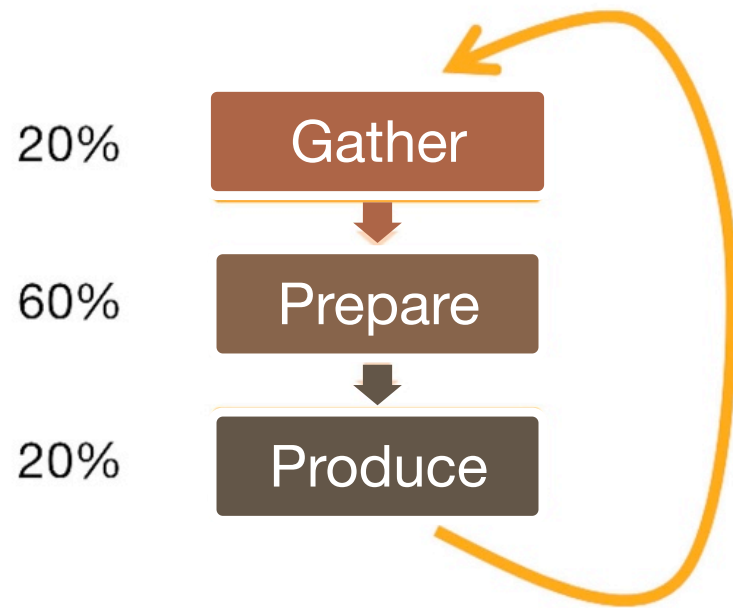# Data percolation: A model of data preparation and analysis



```
Gather
  ↓
Prepare
  ↓
Produce
```

See also: the Data Journalism Handbook

Data percolation: A model of data preparation and analysis

Gather

Prepare

Produce

# How should I budget my time?

20% **Gather**

60% **Prepare**

20% **Produce**

Finding a story is a creative process.

Let it percolate!

SENT DATA

BREAKING NEWS

RECURRING EVENTS

SHARE DATA

THEORIES TO BE EXPLORED

WHAT TO COMPARE OR SHOW CHANGE

WHAT DOES THE DATA MEAN

WHAT OTHER DATA SETS TO USE WITH IT?

SPREAD SHEETS

SPREAD
SHEETS

DATA IN WRONG
FORMAT

MERGED
CELLS

UNNECESSARY
COLUMNS OF DATA

DATA MEASURED IN
DIFFERENT UNITS

PERFORM
CALCULATIONS
ON THE DATA

RECALCULATE
IF NEEDED

SANITY CHECK
THE RESULTS

# How should I budget my time?

| Gather | 1.1 **FIND** reliable data sources |
| --- | --- |
| | 1.2 Understand your **RIGHTS** |
| | 1.3 Visualise and **UNDERSTAND** your data |

| Prepare | 2.1 **CLEAN** your data |
| --- | --- |
| | 2.2 **TRANSFORM** it where useful |
| | 2.3 **COMBINE** it with other data sets |

| Produce | 3.1 **REDUCE** and find the story |
| --- | --- |
| | 3.2 Think and understand the **CONTEXT** |
| | 3.3 Do your results pass a **SENSE-CHECK**? |

# Time planning

Gather

Produce

Prepare

2.1 **CLEAN**

2.2 **TRANSFORM**

2.3 **COMBINE**

2.4 **ENRICH**

2.5 **ANALYSE**

Validating and cleaning data

Inception points: cross filters, outliers and comparisons

# Prepare (Stage 1)

Prepare

2.1 **CLEAN**

2.2 **TRANSFORM**

2.3 **COMBINE**

2.4 **ENRICH**

2.5 **ANALYSE**

# Introducing Open Refine



http://openrefine.org

# Exercise

Validating and cleaning data exercise

Focus on sections:

- 2) Multiple representations

- 4) Summation records

- 5) Mixed use of numerical scales

Thank-you