# Managing Data & Databases

Session 9
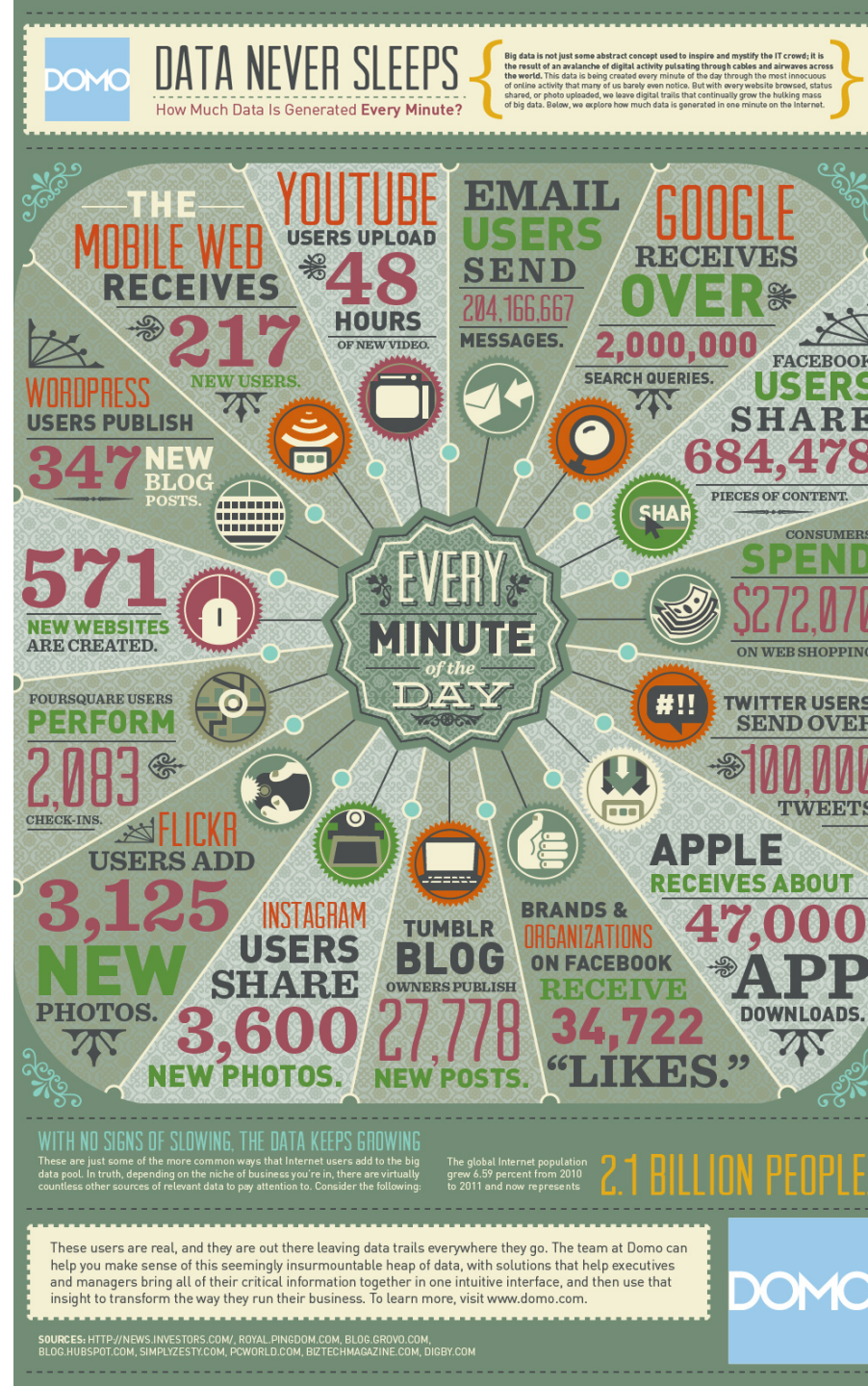The Data Deluge and Big Data

# What is big data?

- The 3 Vs
  - Volume
  - Velocity
  - Variety

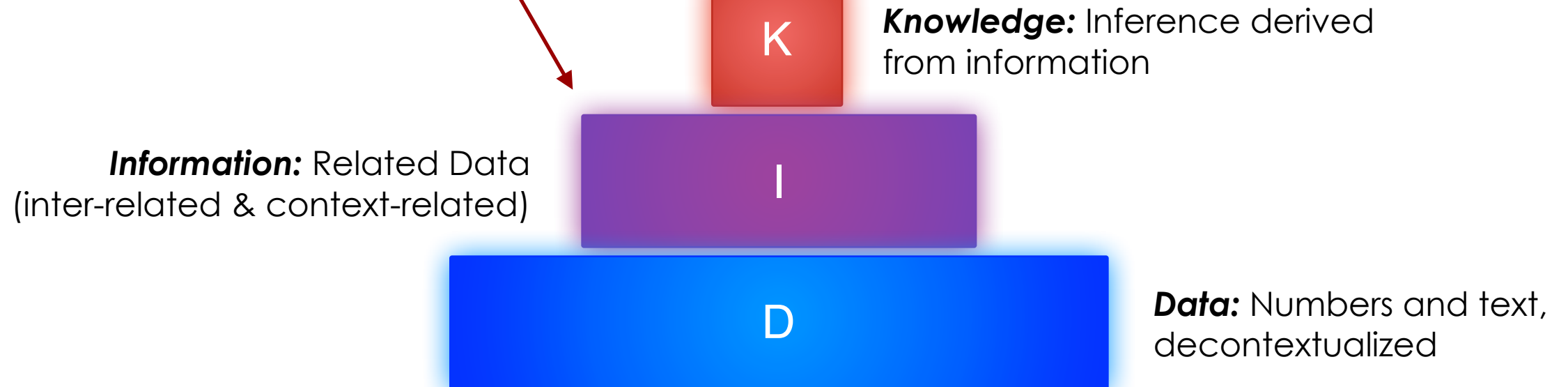- Other characteristics
  - Versatility
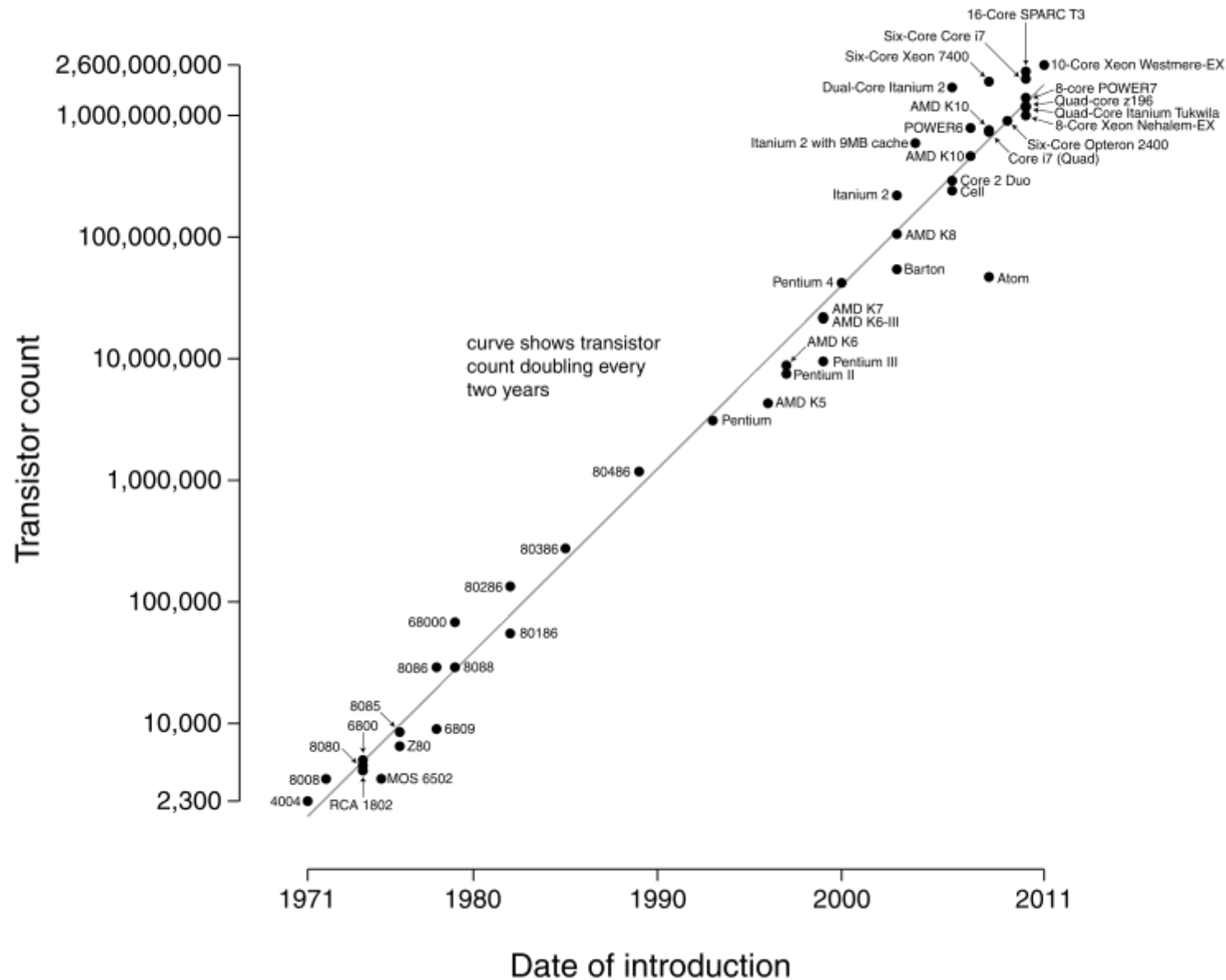  - High Granularity
  - Link-ability

# The 3Vs illustrated

A sneak peek:
http://tweetping.net



Source: domo.com

The bottleneck(s)

K

**Knowledge:** Inference derived from information

I

**Information:** Related Data (inter-related & context-related)

D
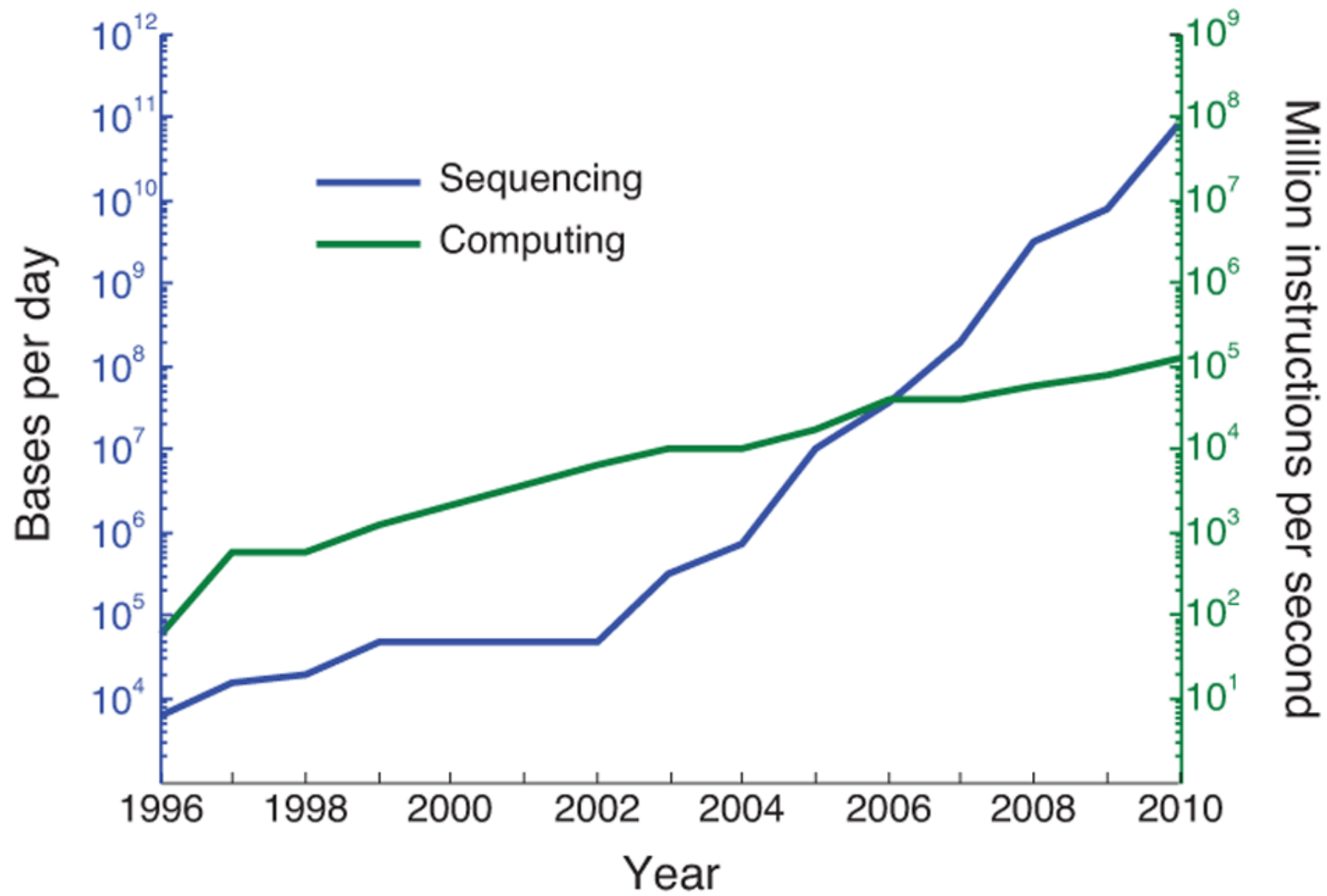
**Data:** Numbers and text, decontextualized

# Microprocessor Transistor Counts 1971-2011 & Moore's Law



Source: Wikipedia

# An estimate
# How much data do we generate?

**2010:** 1.2 zettabytes

**2011:** 1.8 zettabytes

How much is 1,800,000,000,000,000,000,000 bytes?

That fills 115 billion 16GB iPads!

# Where does the data come from?

- Digitization (of existing or new data)

  - Content: Movies, Pictures, Documents, etc...

  - Sensing: DNA, Weather, GPS,... and Apple Health Kit!


- Transaction and Interaction Data

  - Banking, Commerce, Chat,


- User-Generated Content and Digital Footprint

  - Social Media, GPS, Comments, Website visit data,... and Apple Health Kit!

# Are there any perils?

Big data: are we making a big mistake?
By Tim Harford

# How to deal with it?

- Scalable and distributed data structures

- Parallel Computing

- Peer-to-peer Grids

# How to make sense of it? (1)



- Information Design
  - Visualization
  - Infographics
  - Word Clouds
  - etc.

Source: Wikipedia

# How to make sense of it? (2)

- Artificial Intelligence

  - Machine Learning (Pure or Supervised)

  - Pattern Matching

  - Data Mining (Classification, Clustering, etc...)

  - Text Mining (Natural Language Processing)

- http://www.ibm.com/smarterplanet/us/en/ibmwatson/what-is-watson.html

# How to make sense of it? (3)

- Analytics (Adapted Statistics)

  - Large data validity

    - Parametric, Non-parametric, Simulation-Based

    - Inferential vs. Predictive

    - Theoretical vs. A-theoretical

  - Parallel implementation

  - Mixed-level methods

  - Longitudinal methods

# But should we keep all that?

And is there any way not to?