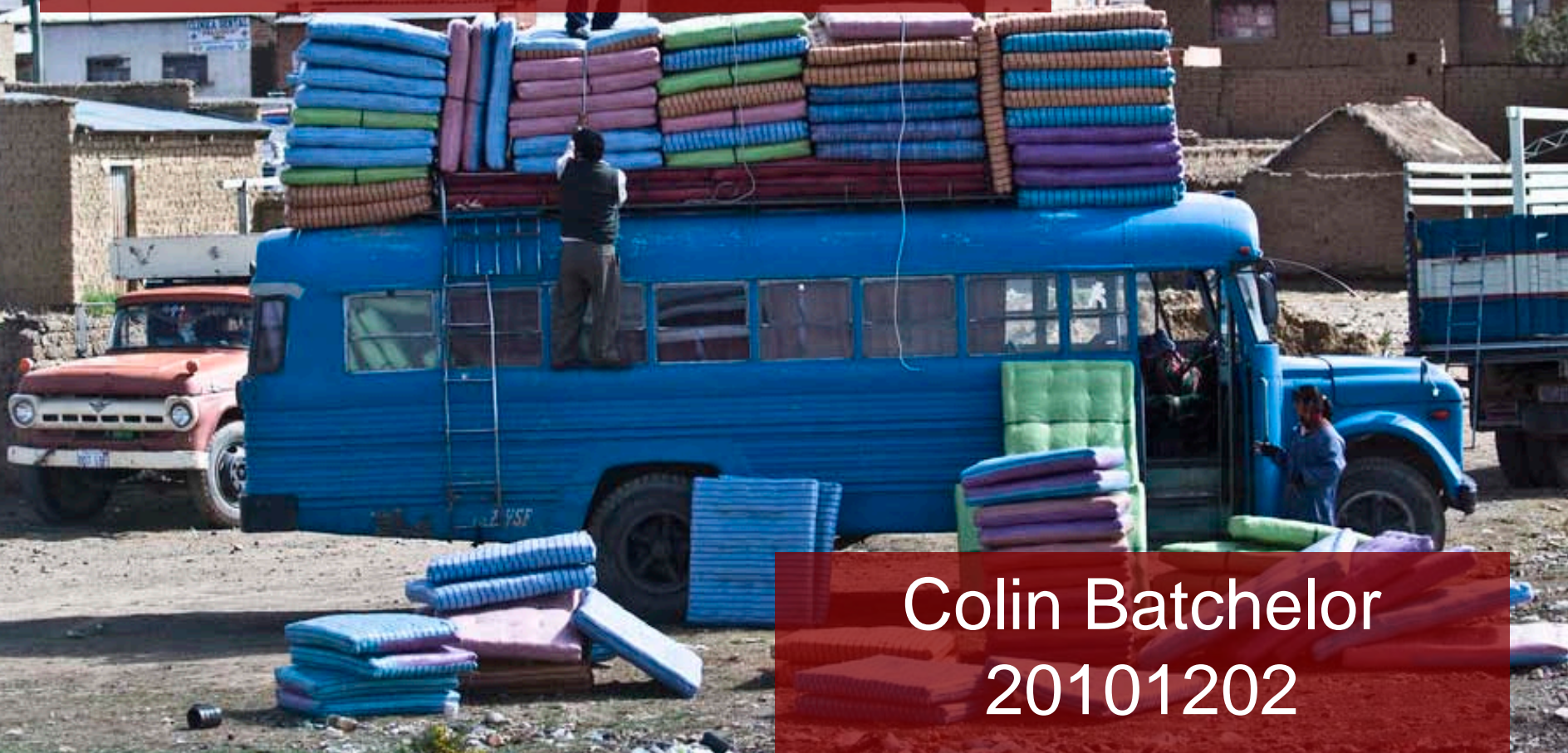


All aboard the semantic  
bandwagon...  
... hold on a moment



Colin Batchelor  
20101202

# What do people mean by semantics?

## Semantic web

- Specifying how web service shall speak unto web service

## Natural-language semantics

- Lexical semantics
- Compositional semantics

# What do search engines do?

**Words** stand only for themselves (Spärck-Jones)



You shall know a **word** by the company it keeps  
(Firth)



Search engines don't use semantic methods.

# Statistics *versus* semantics

A counterblast:

Alon Halevy, Peter Norvig and Fernando Pereira, “**The Unreasonable Effectiveness of Data**”, *Intelligent Systems*, 2009, **24**, 8.

But STM publishers don't just publish text.



Summary of preferred values of parameters for uptake on ice surfaces

# Beyond text

Species	$\alpha/\gamma$	$\pm \Delta \alpha_s$	$K_{inc}$ cm	$N_{max}$ molecule $cm^{-2}$	$\Delta(E_{ads}/R)$ $\Delta \ln N_{max}$	Te
O	$7 \times 10^{-6} + 2.6 \times 10^{-24} \exp(1370/T)$ [O <sub>2</sub> ]	$\pm 0.3$				110-115
O <sub>3</sub>	$< 1 \times 10^{-5}$	$0.7(\Delta \log \gamma)$				220-225
OH			No recommendation			
HO <sub>2</sub>			No recommendation			
H <sub>2</sub> O <sub>2</sub>	0.02	$0.5(\Delta \log \gamma)$	1.6	$3 \times 10^{14}$	$\pm 0.5$ ( $\Delta \log K_{inc}$ )	228-230
H <sub>2</sub> O			No recommendation			
NO	$\leq 5 \times 10^{-6}$	$1.0(\Delta \log \gamma)$				
NO <sub>2</sub>			$3.07 \times 10^{-09} \exp(2646/T)$		$\pm 100$	
NO <sub>3</sub>		$0.5(\Delta \log \gamma)$				170-210
N <sub>2</sub>			No recommendation			190
HONO	0.02	$\pm 0.01$	$1.0 \times 10^{-05} \exp(3843/T)$	$3 \times 10^{14}$	$\pm 50$	180-210
HNO <sub>3</sub>			$1.0 \times 10^{-5} \exp(4585/T)$	$2.7 \times 10^{14}$	$\pm 700$	190-210
HO <sub>2</sub> NO <sub>2</sub>	0.15	$\pm 0.10$	No recommendation			190-210
N <sub>2</sub> O <sub>5</sub>	0.02	$\pm 0.01$				190-210
SO <sub>2</sub>	$0.9 \times 10^{-6}$	$\pm 0.5 \times 10^{-6}$	See data sheet			210-215
HCHO			0.7	$2.7 \times 10^{14}$	$\pm 0.3$	198-200
HCOOH			$5.8 \times 10^{-11} \exp(6500/T)$	$2.2 \times 10^{14}$	$\pm 300$	187-200
CH <sub>3</sub> CHO			No recommendation			
CH <sub>3</sub> COOH			$1.9 \times 10^{-11} \exp(6660/T)$	$2.4 \times 10^{14}$	$\pm 300$ 0.1	195-200

Data which doesn't come in sentences:

- Tables in articles
- Microarray data
- Experimental protocols
- Crystallographic data
- Chemical structures

# The case for chemical markup

There are **very well established and understood** methods for indexing and searching chemical structures.

Chemical structures are **the chief route into the literature** for many chemists.

We can extract chemical structures at **least semi-automatically** from:

**Molecule names**  
**Author-supplied graphics**

# 5 questions for you

1. What could you achieve through “letting words stand for themselves”?
2. What do you hope to achieve through semantic markup?
3. What kinds of non-textual data do you have access to?
4. What can you do with your non-textual data?
5. What external resources could be opened to your readers?

Any questions for me?