

UNIVERSITY OF
BIRMINGHAM

U

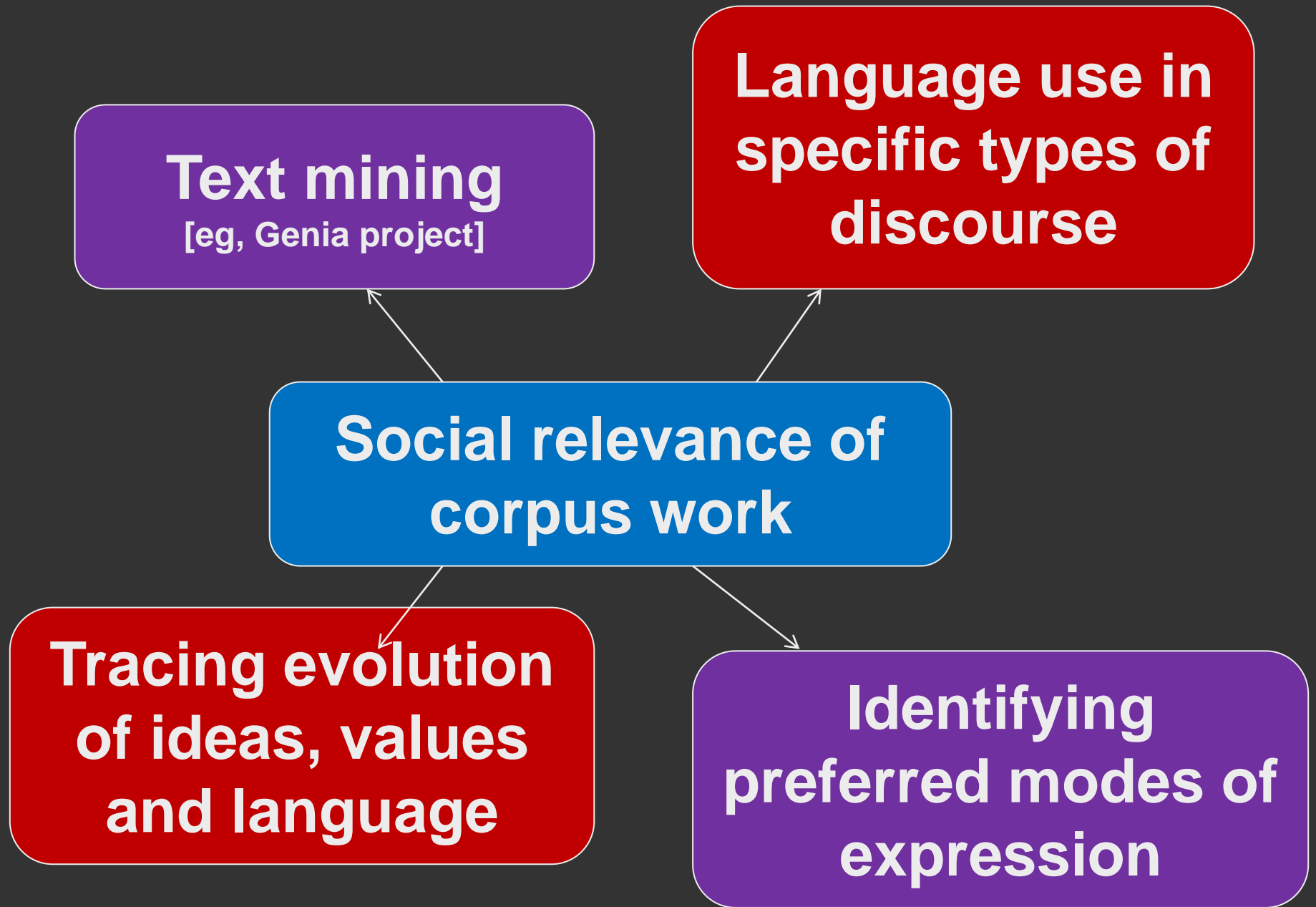
Supporting
Social
Scientists
in the iResearch Age

B

Paul Thompson
Director, Centre for Corpus Research

iUser: [it's not all] About me

- **British academic**
- **Applied linguist**
- **Corpus linguist**
- **Journal editor**



Technology and corpus linguistics



Noam Chomsky:

the linguist must model language *competence* rather than *performance*
No corpus of evidence can represent all language

In the early sixties, with development of mainframes, the first major corpora were compiled

- The Brown corpus was the model: 1 million words, with 500 samples of 2000 words each













Growing up

- The next generation of corpora:
 - COBUILD corpus 1985 = 18 million words
 - Bank of English 1991, increasing to 525 million words in 2005
 - British National Corpus 1994, 100 million words [fixed size]
 - Oxford English Corpus – now 2 billion words



Access

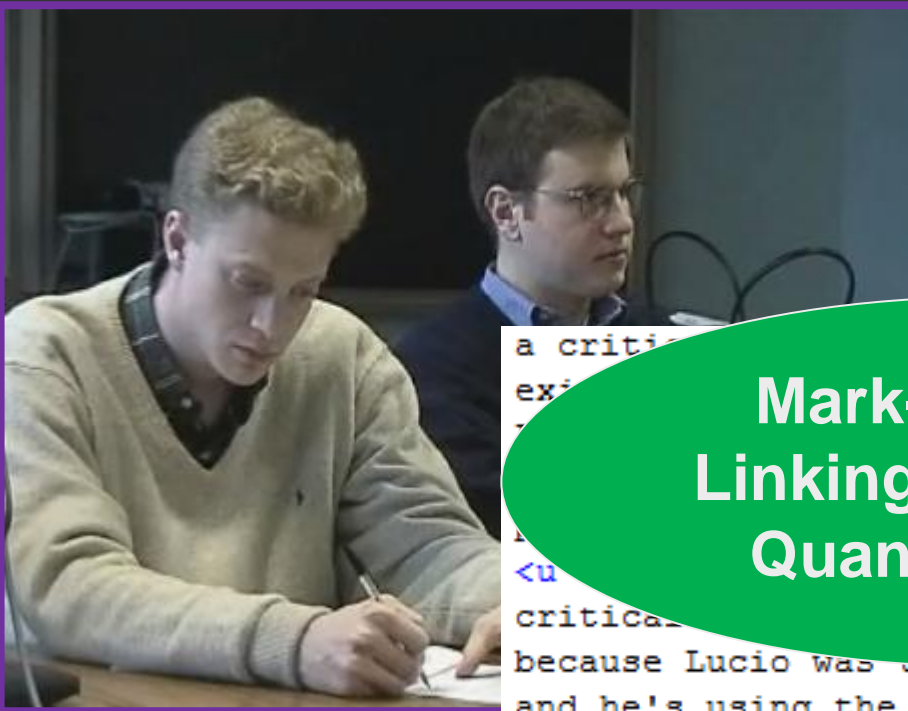
- **In 1960s:**
access highly restricted
- **In 1990s:**
corpus analysis tools limited, but some network access possible;
first ICAME CD-ROM collection of corpora
- **In 2010:**
huge corpora accessed through web interfaces;
many tools for use on own machine;
some owners restrict access to own researchers

Corpus name	Language	Size	
Internet-ZH	Chinese, Simplified	277,931,664	 
British National Corpus	English	112,181,850	 
ukWaC v1.0 old	English	1,526,599,198	 
French web corpus	French	126,850,281	 
deWaC	German	1,627,169,557	 
JpWaC	Japanese	409,384,405	
Russian web corpus	Russian	187,965,822	
Spanish web corpus	Spanish	116,900,060	

Sketch Engine (used by several publishers)

Corpus name	Language	Size	
Arabic web corpus	Arabic	174,239,600	 
Chinese GigaWord 2 Corpus: Mainland, simplified	Chinese, Simplified	250,124,230	 
Chinese GigaWord 2 Corpus: Taiwan, traditional	Chinese, Traditional	455,526,209	 

Next generation: Multimodal corpora



Mark-up of data
Linking of channels
Quantity of data

a critical... the system as it
exi... ing system and
... way the system
...olve the
<u... myself as a
critical... I'll explain why
because Lucio was trying to solve the problem of development
and he's using the existing theories the existing theories
to develop and critical theory is exactly getting out of
those theories right and maybe like developing a new problem
solving theory and that's critical thinking it's not problem
solving</u>
<u who="sm5442"><gap reason="inaudible" extent="1 sec"/>
your kind of # imperialistic notion by</u>

Please feel free to try out some of the searchable online corpora that I've created, including the following ([more information](#)):

[Corpus of Contemporary American English \(COCA\)](#)

410+ million words (US, 1990-2010)

[Corpus of Historical American English \(COHA\)](#) **NEW**

400+ million words (US, 1810s-2000s)

[British National Corpus](#) *

100 million words (UK, 1980s-1993)

[TIME Corpus](#)

100 million words (US, 1923-present)

[Corpus del Español](#)

100 million words (1200s-1900s)

[Corpus do Português](#)

45 million words (1300s-1900s)

Mark Davies:
Brigham Young University

* BYU-BNC: my architecture and interface

Technological developments have led to:

- **New forms, and new levels of access to data and tools**
- **Much larger datasets**
- **... and new approaches to collaboration ...**

Collaboration

- Using **Web 2.0** technologies and values, e.g.:
 - Blogs, clouds, wikis, social bookmarking
 - Open source software
 - Joint construction
 - Increased levels of access for more people
- **Web 3.0** – interconnection of databases, and the Semantic Web

Key Word in Context (KWIC) searches

	Freq
p/n of	163
p/n the	77
p/n in	66
p/n gene	63
p/n .	
p/n ,	
p/n for	
p/n is	
p/n an	
p/n and	26
p/n facial	20

words occurring to left or right of key word

reaction to the change in their mother's **expression** makes it very hard for their mothers to

expressions into equation (2) gives: **expression** may be simpl

transcribed from the phyA promoter this **expression** occurred at

created from expression builder window and an **expression** of (5 × 500 ×

nonmotile strain (FDR875) regulates the **expression** of 2363 gen

using the expression builder again and an **expression** of [Torsiona

taken for granted assumptions. They are the **expression** of a change

flowering time in terms of morphological **expression** of a plant body controlled by particular

regulation (activation or repression) of gene **expression** of a variety of homeobox containing transcript

. Misexpression of Wnt3a induces ectopic **expression** of AER-specific genes, and application

advantage of a second niche for replication. **Expression** of an autophagy-like process facilitates

. The spatial and temporal regulation of **expression** of any given maternal effect gene often

also has problems in recognising facial **expressions** of basic emotions and my find it difficult

expression is also significantly reduced.23 **Expression** of BMP signalling targets are increased

chromosomal abnormalities will prevent the **expression** of certain genes, somatic cell mutations

cotyledons in our CHS:Gus plants, as CHS drives **expression** of chalcone synthase, an enzyme involved

Key word/phrase placed in centre and lines sorted by the first word to the right - this helps the analyst to identify patterns

Pattern searching

complex number systems $\langle p \rangle$ Consider an
percentage elongation (Equation 2) is an
percentage elongation (Equation 2) is an
promoter, we aimed to study the changing
r 355 was selected to allow constitutive
r 355 was selected to allow constitutive
second promoter CAB is normally drives
fact that Gli3 activity is not required for
ing graph is a QQ plot for the first gene
lation (activation or repression) of gene
lation (activation or repression) of gene
sure that when this is the case, greater
dermal development by directly induces

- expression of the form
- expression of the ductility
- expression of the ductility where l_f
- expression of the chlorophyll
- expression of the reporter
- expression of the reporter gene
- expression of the chlorophyll
- expression of the genes
- expression of the control
- expression of a variety
- expression of a variety of homeobox
- expression of the binding protein
- expression of the GATA

$\langle /p \rangle \langle p \rangle$ where all the are arbitrary reals
where l_f is the distance between gauges
is the distance between gauges marks when
a/b binding protein over time after exposure
gene
in a
a/b
de
agai
of homeobox containing transcription factors
containing transcription factors. A homeobox
or in our case the reporter gene luciferase
factors MED-1 and MED-2, which are required

the pattern:
"expression of" + 0-2
words + a noun

Need for support

- **Severe cuts in research council funds expected**
- **Problems of sustainability**
 - **End of funding for AHDS in 2008**
- **Quantitative databases valued**
- **Language databases not valued**
- **There is a greater role for publishers to play**

What roles?

- **Building repositories for data collections**
- **Developing interfaces and analytical tools for researchers, from a range of fields**
 - **Same dataset, different interfaces**
- **Offering increased access to these resources**
- **Facilitating communications between researchers within one field of research, and also researchers in different disciplines**

Why and how?

- Leadership in sustaining and supporting high quality, international research work
- Ensuring sustainability
- Through partnerships with research institutions