

Introducing XAIRA...



An XML-aware tool for corpus
indexing and searching

Lou Burnard

Tony Dodd

Research Technology Services, OUCS



Topics

- ✱ Background: from SARA to XAIRA
- ✱ Architectural issues
- ✱ What can you do that's *fantastic*?

<http://www.oucs.ox.ac.uk/rts/xara/>



Software development: the conventional wisdom

- i. Assess user needs/requirements
- ii. Prototype systems to fit user needs
- iii. Evaluate against user performance
- iv. Repeat from stage ii. until either
 - a) *user is happy, or*
 - b) *money runs out*



Software development: the usual practice

- ☀ Creeping featurism

- ☀ hey, that's a cool idea, I'll bolt that on too

- ☀ The Hausmann effect

- ☀ this is hopeless, we need to drive a few boulevards through here

- ☀ Modularity and standardized interfaces are your only friends



Historical Background (c.1994)

☀ Original design goals

- ✱ robust searching of very large (c. 1 Gb) amount of SGML data
- ✱ re-use available indexing tools
- ✱ usable by researchers in CL, NLP, lexicography

☀ Original assumptions

- ✱ client/server architecture
- ✱ index build once only
- ✱ one specific corpus (the BNC) only



Historical Background (c.2002)

☀ Design goals

- ☀ robust searching of any amount of XML data
- ☀ offload processing to other components wherever possible
- ☀ assume nothing about input DTD

☀ Architecture

- ☀ client/server still valid
- ☀ expect to re-index often
- ☀ expect multiple interfaces



Why another search engine?

☀ Can't you do all this with Google?

- ✱ Digital texts are not just for discovery and display
- ✱ The methods of corpus linguistics have a wider relevance

☀ Can't you do all this with eXist?

- ✱ Probably, but only if you have a team of programmers at your disposal!



Xaira: the key features

- ☀ Supports word search, concordance generation and manipulation, collocation, lexical analysis
- ☀ Uses XML annotation to the max
- ☀ Supports XML-aware complex queries
- ☀ Leverages existing standards
 - ✱ TEI/XCES
 - ✱ Unicode
 - ✱ CSS and XML
 - ✱ SOAP (xmlrpc)
- ☀ Uses efficient and compact indexing appropriate to small or huge corpora



Architectural issues

How do the various parts of a
XAIRA system interact?



First catch your corpus...

- ☀ any collection of well-formed XML documents
 - ✱ if a DTD is supplied, the corpus must be valid
 - ✱ if no TEI header is present, one will be created
- ☀ the more you put in, the more you get out
- ☀ "texts" are defined independently of file structure, as are the relevant units within them
- ☀ all indexing information is stored in the corpus header



Building the indexes

☀ tokenization

- ☀ implicit, following Unicode rules (locale-sensitive)
- ☀ explicit, following mark up
- ☀ supports lexical features (eg collocation)

☀ lemmatization and POS tags

- ☀ special case of "additional key" mechanism
- ☀ generalized to provide fast context-specific searches

☀ tag indexes

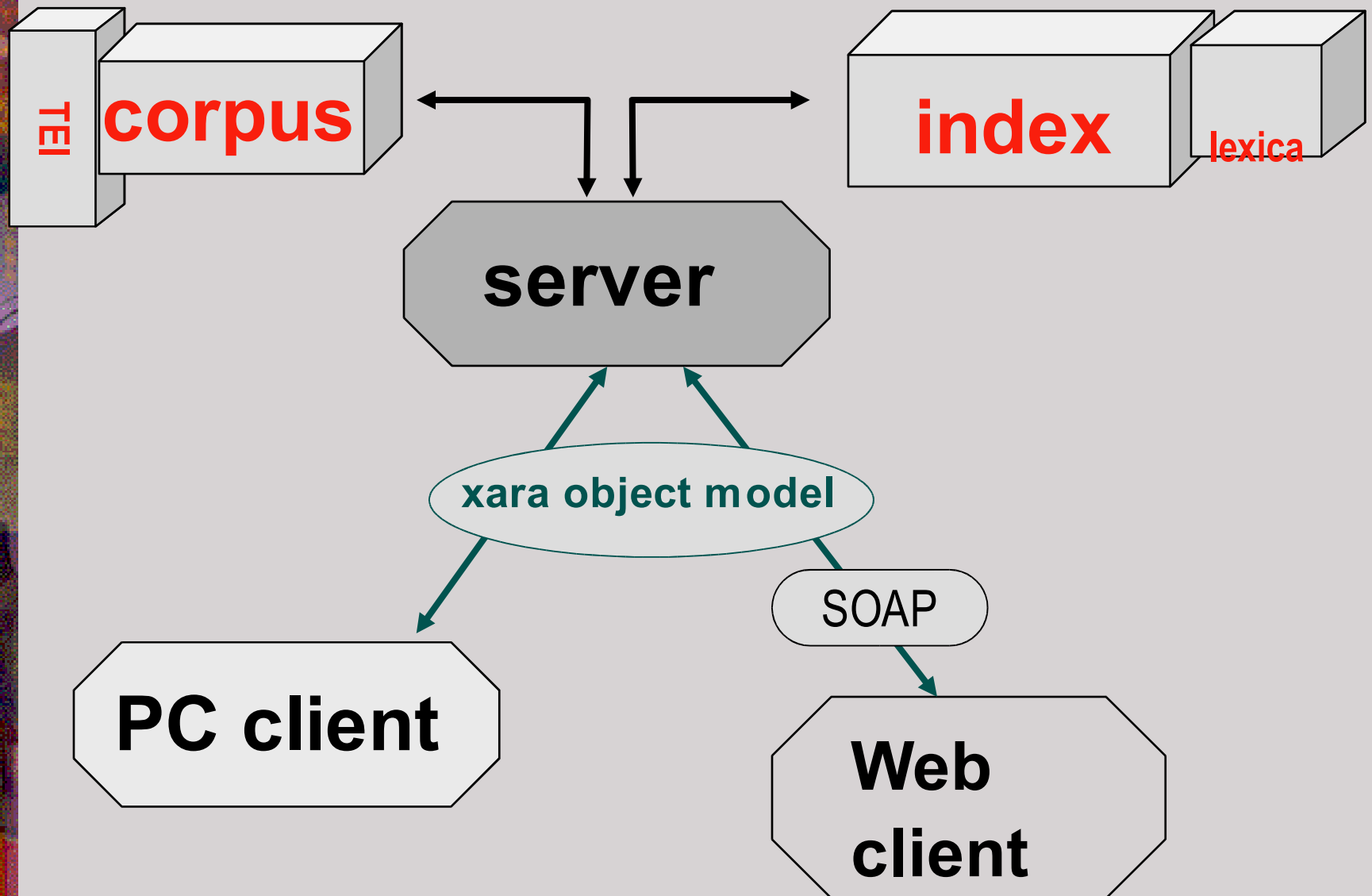
- ☀ attribute values and codebooks



Next, build your index...

- ☀ Can be done simply by adding appropriate declarations to the TEI Header and running the indexer utility
- ☀ But probably easier to do with the supplied *Indextools* utility which
 - ✱ organizes and validates the files you are using
 - ✱ updates (or creates) the header with
 - ✱ *tokenization and indexing rules*
 - ✱ *tag and attribute usage, descriptive codebooks etc.*
 - ✱ *"bibliographic" metadata*
 - ✱ *default behaviour for character encoding, formats used, etc*
 - ✱ optionally runs and tests the indexer

Architecture





Hoorah for Unicode

- ☀ All data is held internally as Unicode
 - ✱ this allows us to defer most problems (e.g. tokenization, case-folding, line-breaking, character normalization, glyph composition) to someone else!
- ☀ User interface issues
 - ✱ For output, use one or more appropriate fonts
 - ✱ For input, we provide a keyboard definition utility



Client/protocols

- ☀ The original SARA protocol
 - ✱ Corpus Query Language
 - ✱ Ad-hoc ASCII strings
- ☀ Now revised completely
 - ✱ Sara Object Model can be accessed
 - ✱ *directly by the client*
 - ✱ *via a SOAP wrapper*
 - ✱ *using saraScript*
 - ✱ The model defines
 - ✱ *CQL in XML*
 - ✱ *methods to manipulate CQL queries and results*
 - ✱ Support for web services



Corpus Query Language

☀ Tokens

- ☀ word, punctuation mark, substring
- ☀ word+annotation/s (e.g. POS)
- ☀ Unicode-compliant regular expressions for words, attribute values
- ☀ XML start- or end-tag, plus attributes

☀ Boolean operations

- ☀ negation, optionality
- ☀ sequence, disjunction, join

☀ Scoped operations

- ☀ within span, within XML element



Client features

- ☀ User-configurable display
 - ☀ plain, XML, user-defined stylesheets
- ☀ User-definable keyboard mapping
- ☀ Texts, Results, Browse windows
- ☀ Results can be exported in XML
- ☀ Scripting language



What can you do that's *fantastic*?

A sketchy over view of Xaira's
query and display facilities



Target queries

- ✱ What is the most frequent noun in this corpus?
- ✱ Find a random sample of 100 instances of "fish" followed by "chips" within 4 words
- ✱ Find sentences beginning with a conjunction.
- ✱ Show all inflected forms of the name "Winston".
- ✱ Show sentences which begin with "well" and end with a question mark.
- ✱ How often and in what contexts is the word "nature" used in different kinds of writing?
- ✱ Which verbs collocate significantly with "bosom" at different periods of history?
- ✱ Do men use colour vocabulary differently from women?



Phrase or simple query

- ☀ search word or phrase
- ☀ can be case sensitive
- ☀ can include punctuation
- ☀ can include *anyword* character
- ☀ watch out for tokenization problems

m mode Criste þenigan æt his halgum weofode, swaswa eowrum hade gerist. Forþan þe ge synd ge-
 wunode on þyssere worulde butan ælcere synne swaswa nan oþer man.
 nað on þam marinum þe næfre wifes ne brucað, swaswa he wunað on þam wifmarinum þe næfre
 and wunedon on clænnysse, Criste fulgigende. swaswa Petrus cwæð to Criste sylfum
 On oþre we sceolan don þe ures drihtnes lare. swaswa he sylf cwæð on his haligan godspelle
 On þam lendenum is, swaswa we leornigað on bocum, seo fule galnys
 and on eallum godnyssum, mannum to bysene, swaswa byrnende leohttatu.
 And Iohannes geseah, swaswa we sædon ær, þone hælend ymbgyrdne æt
 man and modes þæt halige husel him geoffrian, swaswa he sylf getæhte ær his þrowunge.
 era cwellera ehtnyssa, þe þa martyras ofslogan, swaswa we sædon ær.
 num com se haliga gast of heofenum to eorþan swaswa byrnende fyr mid bradum liggette ofer C
 fyr mid bradum liggette ofer Cristes apostolos, swaswa Crist him ær behet
 hyra abhodes; calle hyra þinge him doð gemaene, swaswa him diht se abhod.
 manega Godes þeowan on þam scoton hadum, swaswa us segð se canon.
 am þam mycclan sinoþe, forðan þe se hælend is, swaswa ge gehyrdan oft, soð man and soð God,
 to singanne dæghwamlice urum drihtne to lofe. swaswa se witega Daudi on his witegunge cwæð
 nne, ægðer ge libbendum ge þam forðfarenum; swaswa we leornigað on bocum.
 Man sceal mæssian mid mycelre clænnysse. swaswa þa haligan dydan, þe we hatað confessor
 And mid clænnysse Criste þenedon. swaswa þa canones us cypað openlice Quod nul
 ingð, þæt eow nan syn ne sy, þæt ge swa libban swaswa læwede men.
 ac we mynigað eow, þæt ge clænnysse healdan, swaswa Cristes þegenas on godum gepingðum,
 nas on godum gepingðum, Gode to cwenmysse, swaswa þa haligan dydon, þe we her beforan ræc
 labo eis quod non peribit; God soðlice gecwæð, swaswa us sæde se witega:
 Ac Crist cwæð swa heah swaswa her cwæð on leden. hær synd sume men.



Word Query

- ✱ searches the lexicon for word stem or pattern
- ✱ returns matching word forms with frequencies
- ✱ can restrict by frequency
- ✱ can apply lemmatization rules
- ✱ then carries out a lookup to display hits

so[þð].*

U

Lookup

☒ Return

Word	Frequency	Forms
sop	89	26
sopfæst	7	5
sopfæstnes	8	6
soplic	1	1
soplice	66	2
sopra	1	1
sopre	1	1
sopsagol	1	1
sopsagu	1	1

Form	pos	Frequency
sod	JN	19
sodan	JD	10
sopan	JJD	7
sodan	JJA	6
sode	JN	6
soda	JN	5

9 words

Save

Query

Show

☒ Controls☒ Forms

Download Lemmata

A lemmatisation scheme determines how individual words are grouped under headwords. Select from the options below.

Apply

Headword



XML query

- ✱ searches for XML start- or end-tags (not elements)
- ✱ start-tags optionally qualified by attribute values
- ✱ uses predefined codebooks (value indexes) if available
- ✱ useful in combination with other queries

XML query

The image shows three overlapping windows from an XML query application.

XML Window: Contains a list of XML tags on the right, including 'availability', 'bibliography', 'creationDate', 'distribution', 'encoding', 'extent', 'fileFormat', 'foreign', and 'id'. On the left, there are radio buttons for 'Start' and 'End', and checkboxes for 'Global', 'Pattern', and 'Attributes'.

Attribute selection Window: Shows a list of attributes: 'profileDesc', 'publicationStat', 'refDesc', 'revisionDesc', 'sourceDesc', 'TEI.2', 'teiHeader', 'text', 'textClass', 'title', 'titleStart', and 'w'. Below this list is a section for 'Attributes' with a list containing 'n'. There are also 'Add', 'Remove', and 'Remove All' buttons.

Attribute Menu Window: Displays a table of attribute values and their frequencies. The table has two columns: 'Value' and 'Frequency'.

Value	Frequency
WULF3,231 03	5
WULF3,231 95	1
WULF3,231 96	2
WULF3,231 93	4
WULF3,232 101	3
WULF3,232 103	2
WULF3,267 11	2
WULF3,267 14	2
WULF3,267 2	6
WULF3,267 5	1
WULF3,267 8	3
WULF3,268 17	1
WULF3,268 13	1
WULF3,268 21	2
WULF3,268 24	2
WULF3,268 27	2
WULF3,268 31	2
WULF3,269 43	3
WULF3,269 47	2

At the bottom of the Attribute Menu window are 'OK' and 'Cancel' buttons.



Building complex queries

- ☀ visual interface
- ☀ *scope node* defines where to look
 - ☀ an XML element
 - ☀ by span
- ☀ *query nodes* define what to look for
 - ☀ word, phrase, POS, pattern, XML, or AnyWord
- ☀ *link types* define sequence in which query node targets should occur
 - ☀ next, one-way, two-way

Sentences beginning with conjunctions

Query builder

OK Cancel Query OK

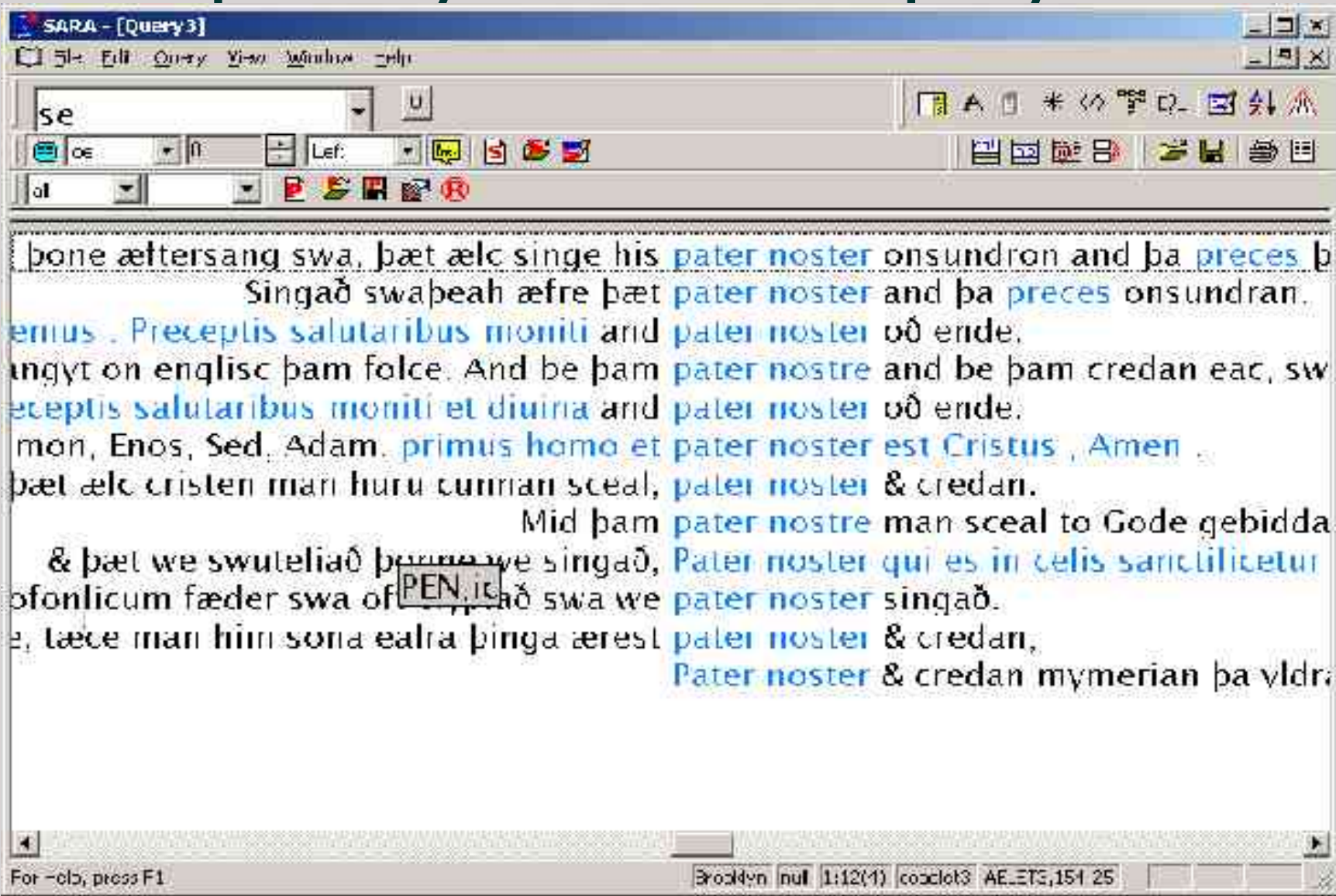
And we ne durran forsuwian, þæt we eow ne secgan, þā deopan lare and yres;
Ac we awengað us,
and þā yfelan mellað.
Ac hi secun swa þeah, þæt him on sumnes sæl lara gesceamige hyra stum,
and mid godum þeawum symble geglencgan
and mid clænum mode Criste þengan æt his halgum weofode, swaswa eowu
and synne on hæpenscype misnotlice gelyfdan
and mid deofles biþgencgum hig sylle forlydan
and þone scyppend forsawon, þe hig gesceop to mannum þurh þæs deofles,
and Moyses hy awrat
and mancynne forbead, þæt hi nænne hæpenscype habban ne mostan, ac sce
And se æ forbead eac swelpe synna
and eac gewitnule þa þe wollice senigdon,
And astealde cristendom
and clænnysse tæhte,
and he sylf is ordfruma eallre clænnysse
and he ana wunode on þyssere worulde butan ælcere synne swaswa nan oþer
And Iohannes se godspeller, þæs heolendes mæg, folode on mæghilde up to
and wunode on clænnysse, Criste folgigende. Swaswa Petrus cwæð to Criste
and we þe folgað.
and we twigan wæran mi þyssere worulde oþþæt Iohannes eom, þe Crist gefullas



Display of results

- ☀ Line (KWIC) or Page mode
- ☀ Context size expandable *ad lib*
- ☀ User defined formatting
 - ☀ stylesheet mechanism based on CSS
- ☀ Export of result files
 - ☀ in XML, or tab delimited

Sample stylesheet display



Collocations

Word Query

god

Look up

Pattern

Word	Frequ...	Forms
god	926	13

Query2

Bvð swapeah sum swa onb

pa godan wyllað

ge synd gesette soðlice to lareow
to lareowum ofer Godes folce, þa
An lyc
e, on þære heahfædera tyman; Of
Ilwæt þa se a
ne to steore sette þa til bec, on þ
ostan, ac sceolde efre wurþian þo
ðlice se soðfæsta hælend, þæs æl
annum þe næfre weres ne brucað,
witega, awrat on his witegunge a
þære ealdan æ. de eac on bære ni

Collocations

Focus

Query: Query2

Hits: 926

Downloads only

<or><lemma>god
</lemma></or>

U

Calculate

Pattern

Node	Frequ...	MI
heorr	2	8.1
wipersaca	4	7.5
þanc	5	7.4
æfestful	1	7.1
marinus	1	7.1
elias	1	7.1
godspellian	1	7.1



Manipulation of results

☀ Sorting

- ☀ by left, right, or centre spans
- ☀ by orthographic form or POS code
- ☀ case sensitive or insensitive

☀ Thinning

- ☀ at random
- ☀ by selection

☀ Analysis and partitioning



Partitions

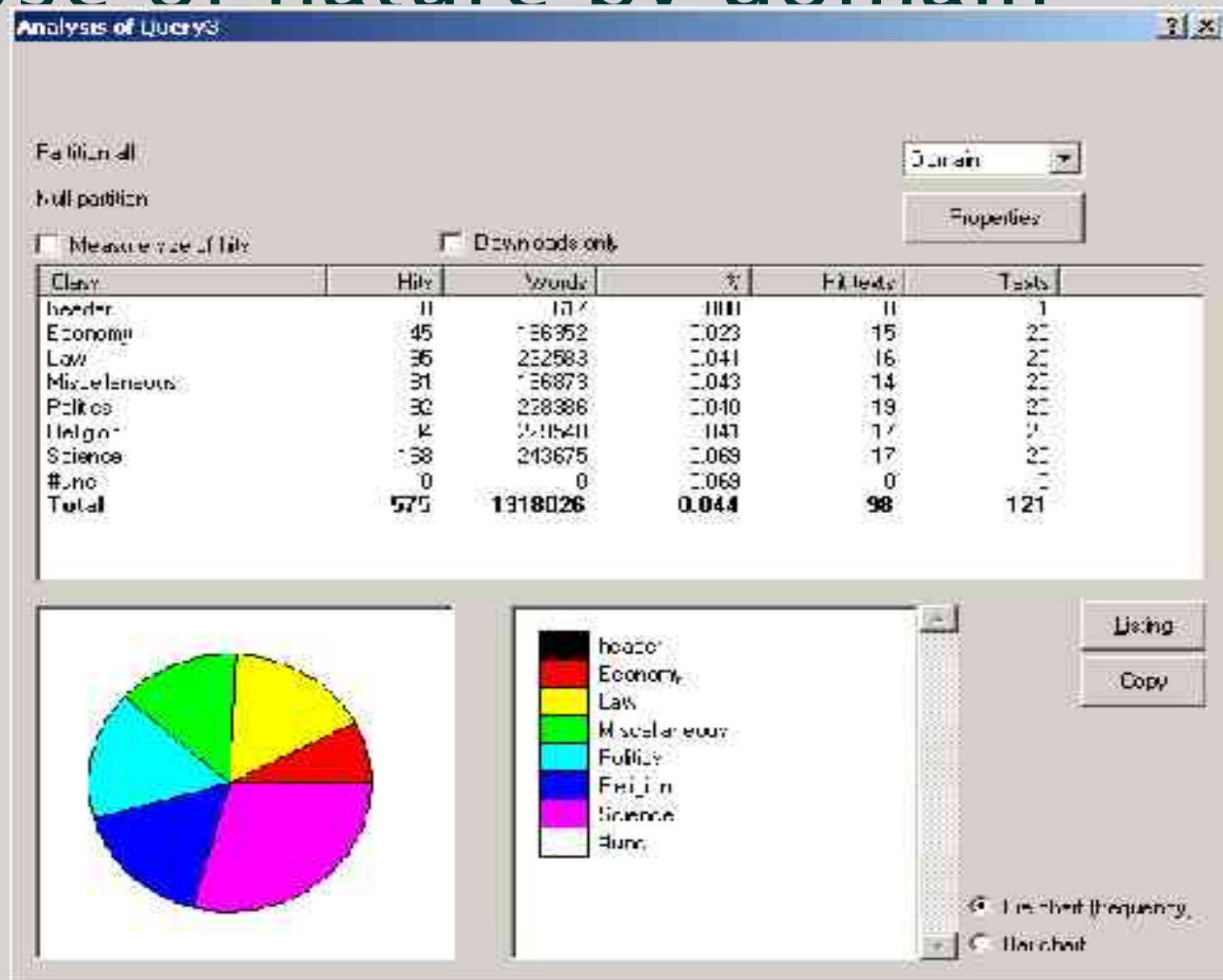
- A partition is a way of grouping the texts making up a corpus, according to
 - *some explicit annotation or characterization (e.g. an attribute value)*
 - *according to whether or not they match a query (a partition of two halves)*
 - *arbitrary manual classification*
- Each member of a partition is a discrete text
- Analysis shows the rate of occurrence of hits within members of the partition
- Partitions can be saved and re-used or defined dynamically
- indextools generates a default partition using `<catRef>` element


crist

Du ner

Text	Class	title	Query
coaelet3	REL	Ælfric's First and Second Letter to Wulfstan	25
cowulf3	HOM	Wulfstan's Homilies (O3)	12
coaelive	BLS	Ælfric's Lives of Saints	8
coaelet4	REL	Ælfric's Letter to Wulfstige	7
cowulf4	HOM	Wulfstan's Homilies (O3/4)	6
colaw2	LAW	Alfred's Introduction to Laws, Alfred's Laws, Ine's Laws	5
cochroa2	HIS	Anglo-Saxon Chronicles until 946	5
colaw3	LAW	Eleventh Century Laws	5
coorosiu	HIS	Alfred's Orosius	1
cogregd4	BLS	Dialogues of Gregory the Great (MS C)	1
coboeth	PHI	Alfred's Boethius	0
cochroa3	HIS	Anglo-Saxon Chronicles 951 - 1001	0
coapollo	FIC	Apollonius of Tyre	0
cobede	HIS	Bede's Ecclesiastical History	0
header	header	Brooklyn Corpus	0
cogregd3	BLS	Dialogues of Gregory the Great (MS H)	0
colaw4	LAW	Late Laws	0

Use of nature by domain





Saving and re-using queries

- ✱ Bookmarks
- ✱ Queries are saved with thinning information
- ✱ Optional annotation
- ✱ Associated bookmarks are preserved