# Visualising Patterns in Text

Mick O'Donnell

Universidad Autónoma
de Madrid

# Thanks…

Thanks to:

Roberto Therón

and

Michelle Zappavigna

...for discussions on this topic

# Note...

Most of the visualisations available here will be available in the next version of UAM CorpusTool

http://www.wagsoft.com/CorpusTool/

(Free, Windows and MacOSX)

# The TREACLE Project

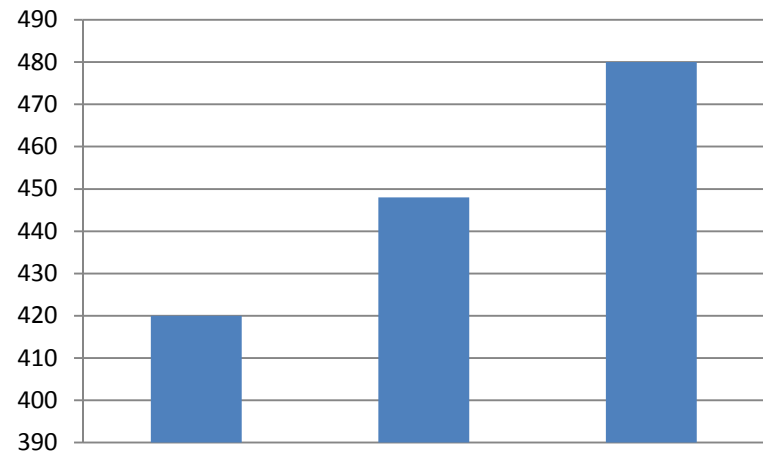- Project: TREACLE

  **T**eaching
  **R**esource
  **E**xtraction from an
  **A**nnotated
  **C**orpus of
  **L**earner
  **E**nglish

  *Official Title: "Developing an annotated corpus of learner English for pedagogical application"*

- A cooperation between:

  Universidad Autónoma de Madrid and

  Universitat Politécnica de Valencia

- Funded by Ministerio de Ciencia e Innovación (FFI2009-14436/FILO)
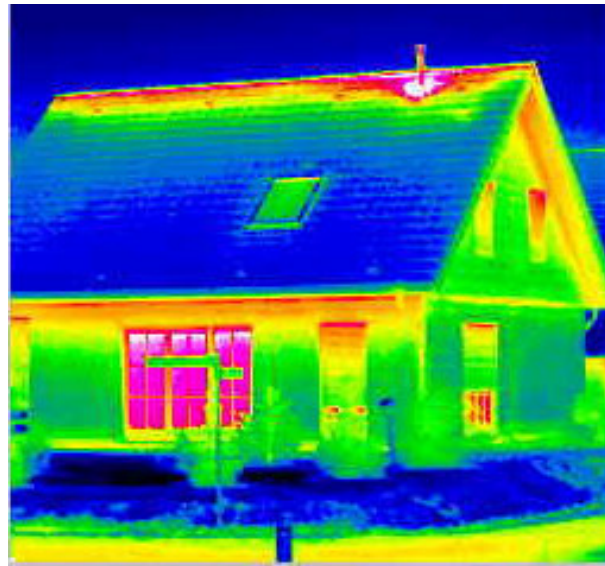
- Runs: January 2010 – December 2012

# Aims of this talk

- Humans can absorb data much more rapidly when it is presented as an image of some kind (e.g., as a bar chart) in comparison with tables of numbers.

# Aims of this talk

- Language is highly structured, but organically complex.

- Appropriate visualisation of linguistic data allows us to see patterns which are otherwise obscured.

- This talk will explore various ways of presenting the patterns in a text in graphical form.

**What can we visualise?**

- In Corpus Linguistics, two main objects to explore:

Patterns over the corpus: patterns observable as generalisations over a large corpus of texts.
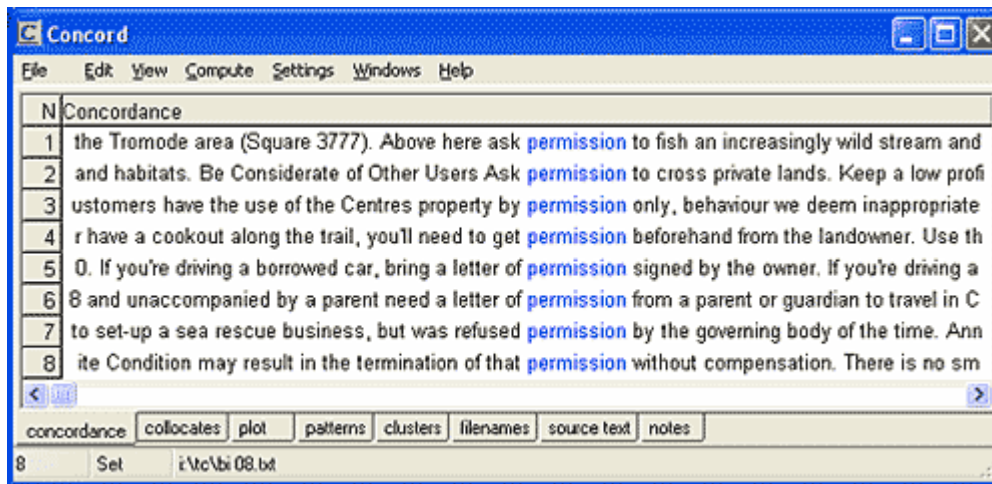
Patterns in a single text: patterns observable by close analysis of an individual text.

# 1. Visualising a Corpus vs. Visualising a Text

## Visualising Patterns in a corpus

**E.g., KWIC tables to show contexts of a lexical pattern**

**E.g., Sketch Engine grammar summaries**

- More Graphical visualisations of KWIC information?



LInfoVis "End To End" shows alternative word sequences between two words over a corpus

http://www.eurac.edu/en/research/institutes/multilingualism/Projects/LInfoVis/default.html )

# 1. Visualising a Corpus vs. Visualising a Text

- However, this talk will focus on visualising patterns in single texts.

- Question behind this talk:

How can we best visualise the data in a text to reveal the underlying patterns within the text?

# 2. What to visualise and How to visualise

- The **variety of information** that can be visualized in a text is quite open, depending on what aspects that one wants to explore

- The manner of **presentation** (visualisation) of any of type of information is also quite open.

- However, some types of information have better affinity to certain visualizations.

- Thus, two questions:

  – **What to visualise?** (what data best shows the patterns in a text?)

  – **How to visualise?** (which means of presenting a particular type of data best reveals the patterns).

# 2. What to visualise and How to visualise

- Visualisations can show:
  - What information is linguistically **important** in some way (wordclouds, highlighted text)
  - How data is **grouped** (e.g., scatter plot)
  - Which entities are **outliers** (e.g., scatter plot)
  - How data is **structured** (e.g., semantic nets, syntactic trees)
  - How data **differs** (e.g., bar charts, pie charts, etc.)
  - How data **differs over time** (e.g., bar charts, flow diagrams)
- So, before we decide how we present out data, we need to decide what we want to reveal (groupings, differences, structure, etc.)

the cat sat

# 2. What to visualise and How to visualise



Visualising text groupings using PCA:

A: Abstract
E: Editorial
N: News Article
F: Fiction

Information used to place texts: degree of usage of syntactic features

Visualising text groupings using PCA:

A: Abstract
E: Editorial
N: News Article
F: Fiction

Information used to place texts: degree of usage of syntactic features

# 3. Visualising Patterns in a Text: Lexical Patterns

How do we explore the LEXICAL patterns in a text?

1. What to visualise?

- Which words are most important in the text? (and what does 'importance' mean?)

- How does the lexis in the text differ from other texts? (keywords)

- Which sequences of words in the text are common artifacts of the field or genre? (phrases, n-grams)

- How does the lexis in the text change as one moves through the text?

# 3. Visualising Patterns in a Text: Lexical Patterns

## 2. How to Visualise (i)

Table Views: provide precise numbers but the data is not so readily absorbed. E.g., a table of "keywords" in a text, sorted by keyness.

| Token | N (Text) | N (Ref. Corpus) | Propensity |
|---|---|---|---|
| palestinian | 5 | 19 | 67.14 |
| gaza | 3 | 7 | 54.67 |
| bus | 3 | 7 | 54.67 |
| map | 3 | 10 | 49.75 |
| jerusalem | 3 | 11 | 48.71 |
| violence | 4 | 22 | 46.39 |
| road | 3 | 17 | 45.02 |
| israeli | 4 | 26 | 39.25 |
| bomber | 2 | 8 | 31.89 |
| prime | 2 | 11 | 30.15 |
| himself | 2 | 15 | 28.92 |
| bombing | 2 | 19 | 26.86 |
| israel | 2 | 19 | 26.86 |
| minister | 2 | 24 | 21.26 |
| attack | 3 | 43 | 17.80 |
| mr | 2 | 32 | 15.95 |
| plan | 2 | 33 | 15.46 |
| killed | 3 | 52 | 14.72 |

In better visualualisations:

.

Properties of the text

Realised via

Properties on the screen

In better visualualisations:  E.g.,

- **Properties of the text**
  **Keyness**

Realised
via

**Properties on the
screen**
**Visual Salience**

# 3. Visualising Patterns in a Text: Lexical Patterns

**Wordclouds**: A common means of visualizing the **important** lexis in a text.

- **Importance** of a word in a text is realised in terms of its **visual salience**.

- Importance measured in terms of:

  – Frequency of the word in the text, ignoring closed class words (Wordle), or

  – Keywordiness: Relative frequency compared to relative frequency in a reference corpus

give word tweak images share a fonts appear like greater create different however use provide layouts text friends clouds save Wordle

# 3. Visualising Patterns in a Text: Lexical Patterns

Visual Salience realised in terms of:

- Size: size of the displayed word relates to the importance of the word.

- Centrality: More important words can appear at the centre, less important on the periphery.

- Hue: e.g., more important words realised by "warmer" colours

- Lightness: less important words fade out.

# 3. Visualising Patterns in a Text: Lexical Patterns

**Wordclouds:** pros and cons

- Wordclouds allow you to quickly perceive which words are important in a text.

- However, they remove words from their context, so it is difficult to see how the word is functioning in the text, what other meanings it relates to, etc.

- A KWIC table can show the contexts of a word, but only one keyword at a time.

- Solution: apply the wordcloud concept  (importance = visual prominence) to the display of the text itself….

➡ **Text-Salience View**

# Text-Salience View: Keyness of words

# Using "subjectivity" rather than "keywords"

Strong-Negative  Weak-Negative Neutral Weak-Positive Strong-Positive

Public opinion surveys report little support in the United States for those who oppose military strikes against terrorists. But protesters' arguments deserve to be taken seriously nonetheless, and not just because they are seriously offered. Many of us would wish that the pacifists were right - that this problem could be solved without more killing - and such wishfulness is likely to translate into more skepticism about a war on terrorism as that war gets tougher.

The United States is entitled to defend itself; the United States is morally obliged to defend itself. That has to be the starting point.

The Sept. 11 attacks may not have originated with one nation, but they were an attack against the American nation, part of a larger effort to cripple the country and its way of life. The attacks were planned and perpetrated by enemies of the United States, and the right response is to attack those enemies and seek to eliminate the threat they pose. Against this view several arguments are offered. One is that U.S. attacks will cause suffering to many innocent people. This is likely true, as it was in the Civil War, World War II and every other conflict; that is why war is always a last resort.

## Text-Salience View

**Problem with Text-Salience View**

- The page is not the best shape for visualising the patterns in a text: adjacent words in reading sequence may be displayed on opposite sides of a page:

Donald Trump has been stealing attention from campaign issues and candidates with his focus on President Barack Obama's birth certificate and school grades, Republican Senators Lindsey Graham and John McCain said on network talk shows today.

"There's a lot of things Mister Trump can be proud of, but some of this rhetoric and this focusing on the president's birth, I do not think is the way for us to win the White House," Graham, a South Carolina Republican, said on "Fox News Sunday."

Visual perception of patterns may thus be disrupted.

## Text-Salience View

- The design team that designed our writing system obviously had to compromise...

The J-Walk Blog is written by John Walkenbach, and is coming to you from Tucson, Arizona. John is a computer book author, an occasional consultant, and is responsible for several popular Excel add-in products

An alternative to spiral texts involves making a U-Turn at the end of each line and continuing back to the left hand-side, except upside-down. This text is an example of such an approach. Would you like to try reading a book in this manner?

**Visualising change through the text**

- One solution is to separate the jobs of displaying sequence patterns in text from the job of displaying the text itself.

# 3. Visualising Patterns in a Text: Lexical Patterns
## Dispersion plots

- Systemic Coder (2001) allowed viewing of linguistic tags through a text (aka Scott's comment on dispersion plots)

## Heat diagrams

- Heat diagrams show occurences of a word thoughout a text

- Warm colors indicate higher frequency, cold colors indicate periods of low or zero occurence.

- Mike Scott (at AELINCO 2011) discussed problem of displaying variation of occurence of keywords throughout a text.

- Last night, he demonstrated dispersion table to show where individuals are mentioned through a novel.

- This is an alternative means of achieving the same end.

- Heat Diagrams can be scaled to fit on the visible screen.

spatial

Start of text

Low freq. occurences

Higher freq. occurences

Endof text

## Heat diagrams

## Summary of Lexical Visualisation

- We have presented several means of visualising lexical patterns in a text.

  - Tables of keywords provide precise information, but are not immediately digestible.

  - Wordclouds provide a quick synopsis of the important words in the particular text but hide their context.

  - Text-Salience views show which words are important AND the context in which the word appears.

  - Heat diagrams show how particular words change in importance as the text unfolds.

# 4. Visualising Patterns in a Text: **Syntactic Patterns**

- Lexical studies have dominated Corpus Linguistics, mainly because software has been able to intelligently handle words without too much trouble.

- In recent years, syntactic parsing has become robust enough to be used for corpus linguistics, without need for human post-editing.

- UAM CorpusTool uses the Stanford parser (Klein and Manning 2003) to parse each sentence in a user's (English) corpus.

- This information can be used to explore other patterns in the text.

- "Key Features" are the syntactic equivalent of keywords: the features in a text or subcorpus which occur relatively more frequently than in a reference corpus.

- We can use some of the same visualisations as for keywords.

| Feature | N (Text) | N (Ref. Corpus) | Propensity |
|---|---|---|---|
| simple-future | 2 | 162 | 29.23 |
| future-clause | 2 | 165 | 28.70 |
| with-connector | 4 | 590 | 16.05 |
| infinitive-clause | 7 | 1033 | 16.04 |
| modal-clause | 4 | 745 | 12.71 |
| relative-clause | 3 | 629 | 11.29 |
| monotransitive-clause | 13 | 2906 | 10.59 |
| nonfinite-clause | 9 | 2154 | 9.89 |
| transitive-clause | 16 | 4060 | 9.33 |
| true-modal-clause | 2 | 573 | 8.26 |
| simple-present | 7 | 2365 | 7.01 |
| relational-clause | 5 | 1732 | 6.83 |
| doing-clause | 16 | 5602 | 6.76 |
| past-participle-clause | 1 | 368 | 6.43 |
| present-clause | 7 | 2643 | 6.27 |
| passive-clause | 2 | 761 | 6.22 |
| not-perfect-aspect | 24 | 9566 | 5.94 |
| ditransitive-clause | 3 | 1210 | 5.87 |

…show which syntactic features this text makes more than normal use of

future-clause

simple-future

with-connector

infinitive-clause

modal-clause

monotransitive-clause

transitive-clause

simple-present

relational-clause

ditransitive-clause

not-progressive-aspect

present-clause

past-participle-clause

positive-clause

not-perfect-aspect

doing-clause

true-modal-clause

passive-clause

relative-clause

nonfinite-clause

- To see how syntactic choices vary throughout a text, we can highlight segments of text which have nominated features.

Since the first of January 2006, smoking in public places, Duch as pubs, restaurants and offices, is forbidden; this is what the new antitobacco law establishes. S, which was the introductory of tobacco in Europe, regarding the antitobacco law, has become one of the most restrictive countries together with Ireland, Norway and Italy. This law, exaggerated for some people and fair for others, has create a very controversial debate that confronts smokers with non-smokers. In this essay, I intend to present different points of view about the new antitobacco law.

This law establishes smoking zones in pubs, restaurants etc. It limits publicity refering to tobacco and hardens the normative of smoking in public places. In addition, it attempts to improve spanish citizens health, as it is a fact that the first cause of death in our country is tobacco. A recent study indicates that 38.5 % of the population agree with this law whereas 3.11 % are aginst it. According to this results, people should considerate that 25.8 % of people smoke, 26.7 % have given up smoking and 47.5 % do not smoke.
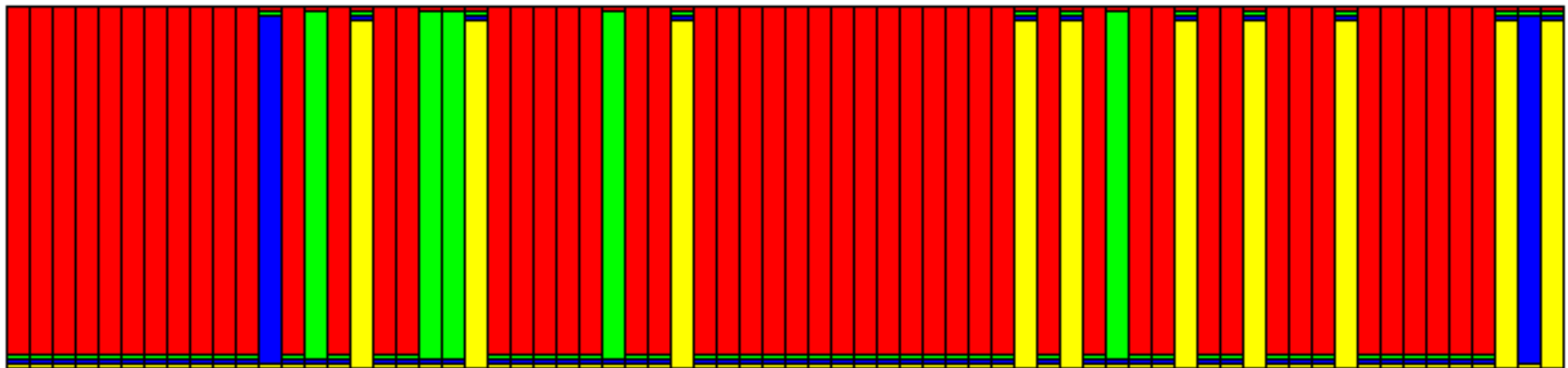
Non-smokers, who are in favour of the antitobacco law, support that the law is going to improve society's health and is not against nobody's rights, in fact, it protects the right to health, which is reflected in the Constitution. Those who do not smoke also support that <Q>. Moreover, they answer to the representatives of the inkeeper sector, who believe that this new law is going to decrease the market share, that this is not going to happen due to the fact that the 70 % of spanish population is non-smoker.

On the other hand, the spokesman of the "smoking club" criticises the Constitution by saying that "the law forces 6 million workers who smoke to go out for consuming a legal product that is being sold by the state. Furthermore, they consider that the law has created a situation which is not fair for smoking people. In addition, they do not understand the government's legislation because it seems that the state wants to reduce tobacco comsumption although it continues selling it and earning money with that.
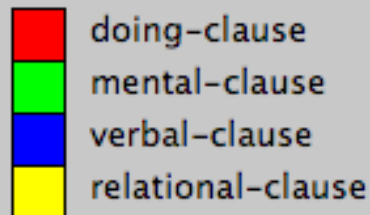
# 4. Visualising Patterns in a Text: **Syntactic Patterns**
Patterns in the unfolding of text

Each bar represents one clause through the text, showing the syntactic feature selected
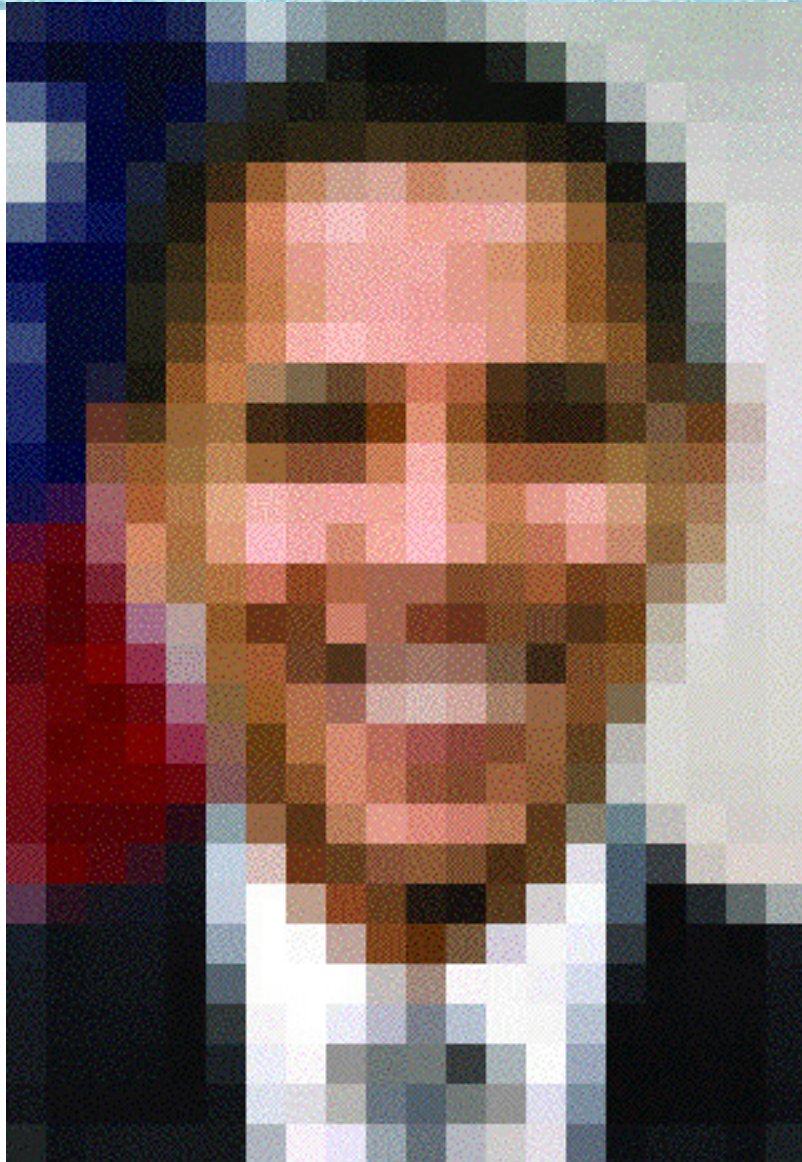


System to graph: PROCESS-TYPE  Smoothing: 0

- 🟥 doing-clause
- 🟩 mental-clause
- 🟦 verbal-clause
- 🟨 relational-clause

## The Value of Smoothing

# The Value of Smoothing

# 1. Visualising Patterns in a Text: Lexical Patterns
## TextFlow diagrams

**Unsmoothed** (each bar represents one unit of text)



System to graph: PROCESS-TYPE   Smoothing: 0

- doing–clause
- mental–clause
- verbal–clause
- relational–clause

**Smoothed** (the value of each segment shared with nearest neighbors)

# Flow diagram plus color-coded text gives best of both worlds



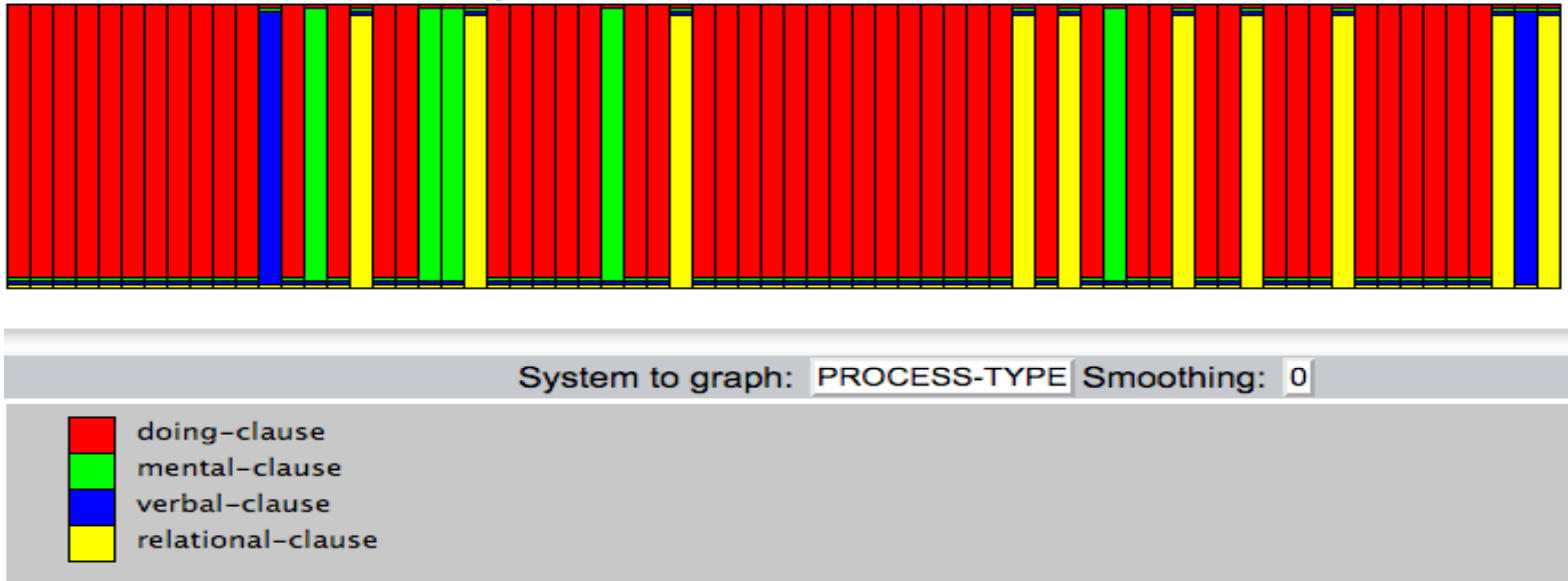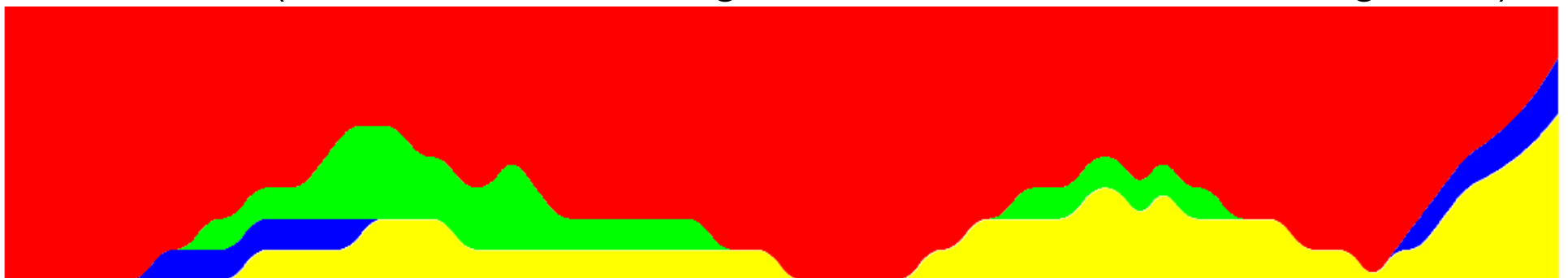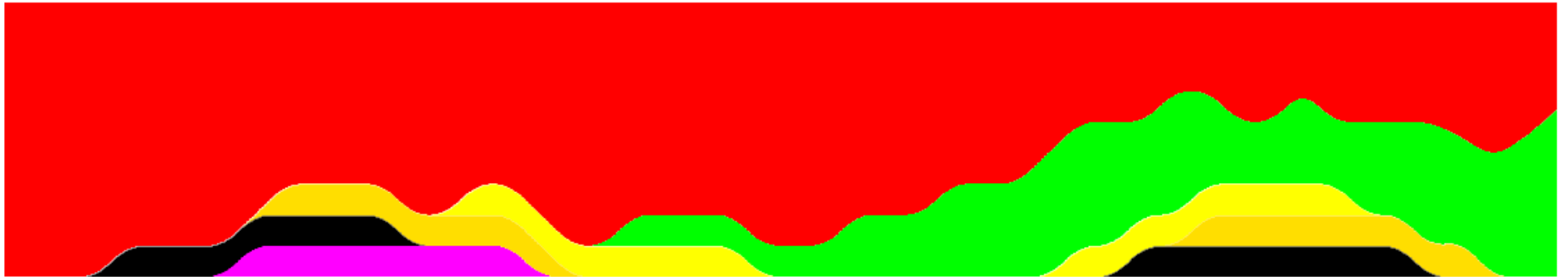KEY: simple-past  simple-present  simple-future  past-perfect  present-perfect  future-perfect  past-progressive  present-progressive  future-progressive
simple-modal  modal-progressive  modal-perfect  past-progressive-perfect  present-progressive-perfect  future-progressive-perfect  modal-progressive-perfect

Moments after Abbott jogged out the back of the Chinook into the rocky, snow-patched valley at about 6 a.m., the yell burst out: "Incoming!"

Enemy rifle fire - sporadic at first, but rapidly building - shot down on the US troops from a fortified ridgeline partway up the mountainside on the east. Rounds also flew in from two dozen black-uniformed Al Qaeda, who appeared to the men like ants, climbing the ridge to the west.

Abbott and the rest of the 10th Mountain Division's Charlie Company took cover in shallow ditches and rises and began returning fire in both directions. To move faster, many dropped their 85-pound rucksacks stuffed with food, cold-weather gear, and extra ammunition - a decision they would later regret.

Battalion commander Lt. Col. Paul LaCamera, the 10th Mountain Division's senior officer on the ground and his top enlisted man, Command Sgt. Maj. Frank Grippe, took up position in the "bowl." The dip in the valley was so fortuitous that Sergeant Major Grippe later joked that Mother Nature must have put it there "because sometime young Americans were going to need it."
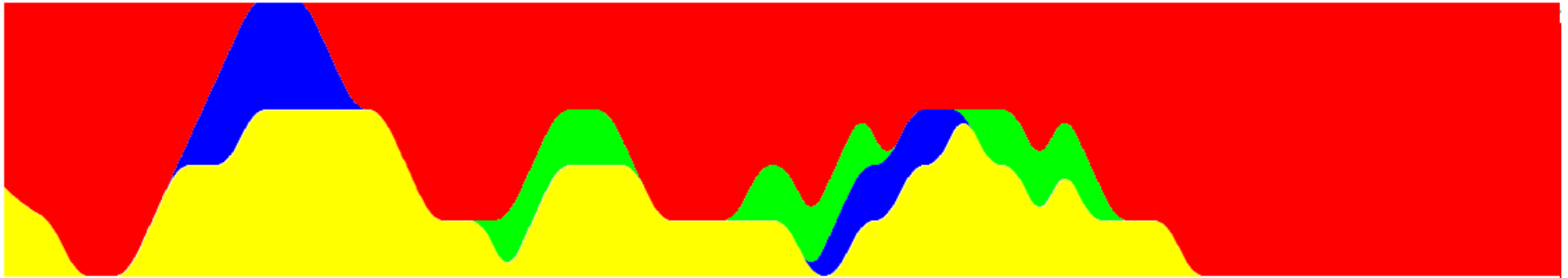
Still, Grippe realized his men were "in a very precarious position." Originally assigned to block major southern and eastern exit routes from the valley - code-named Heather and Ginger - they now faced a "nose-to-nose" battle with a large enemy contingent. Worse, the enemy manned well-concealed mountainside positions, while Grippe's men occupied the exposed low ground.

The consolation prize: They had found one of the biggest concentrations of Al Qaeda and Taliban in the Shah-e Kot Valley.

A few feet from Grippe, Capt. Scott Taylor, the battalion fire-support officer, began urgently radioing for Apache attack helicopters. The bespectacled young officer from Long Valley, N.J., was a 1996 Rutgers graduate in international environmental studies. Now he dealt with the landscape of war, directing artillery and close airstrikes.

Captain Taylor passed grid coordinates for targets to two Apache pilots, who within 20 minutes swooped in firing cannons on the ridge. They circled and were launching rockets at the hilltop when enemy bullets strafed the aircraft, forcing them to leave. An hour later, two more Apaches flew in. Taylor again directed them to targets, but before they fired a single shot, they, too, were repelled by ground-to-air missiles and small arms.
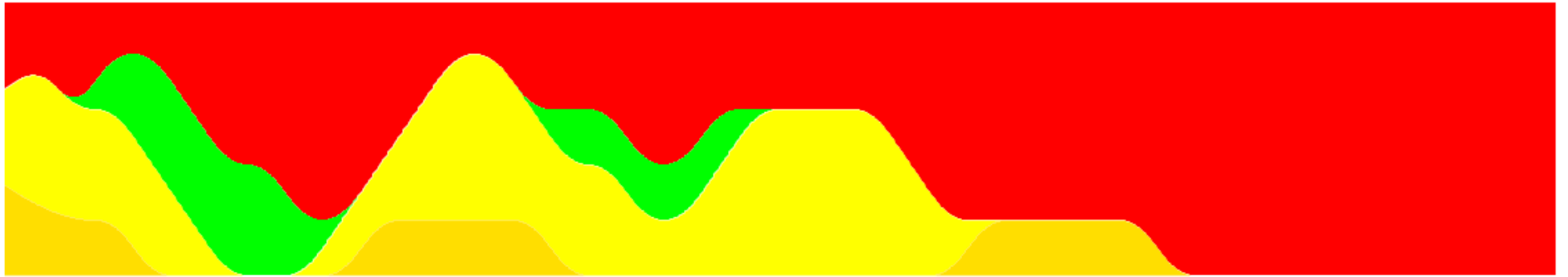
# Oscar Wild The young king

It was the night before the day fixed for his coronation, and the young King was sitting alone in his beautiful chamber. His courtiers had all taken their leave of him, bowing their heads to the ground, according to the ceremonious usage of the day, and had retired to the Great Hall of the Palace, to receive a few last lessons from the Professor of Etiquette. Some of them who had still quite natural manners, which in a courtier is, I need hardly say, a very grave offence.

The lad - for he was only a lad, being but sixteen years of age - was not sorry at their departure. He had flung himself back with a deep sigh of relief on the soft cushions of his embroidered couch, lying there, wild-eyed and open-mouthed, like a brown woodland Faun.

And, indeed, it was the hunters who had found him, coming upon him almost by chance as, bare-limbed and pipe in hand, he was following the flock of the poor goatherd who had brought him up, and whose son he had always fancied himself to be. The child of the old King's only daughter by a secret marriage with one much beneath her in station - a stranger, some said, who, by the wonderful magic of his lute-playing, had made the young Princess love him; while others spoke of an

And, indeed, it was the hunters who had found him, coming upon him almost by chance as, bare-limbed and pipe in hand, he was following the flock of the poor goatherd who had brought him up, and whose son he had always fancied himself to be. The child of the old King's only daughter by a secret marriage with one much beneath her in station - a stranger, some said, who, by the wonderful magic of his lute-playing,
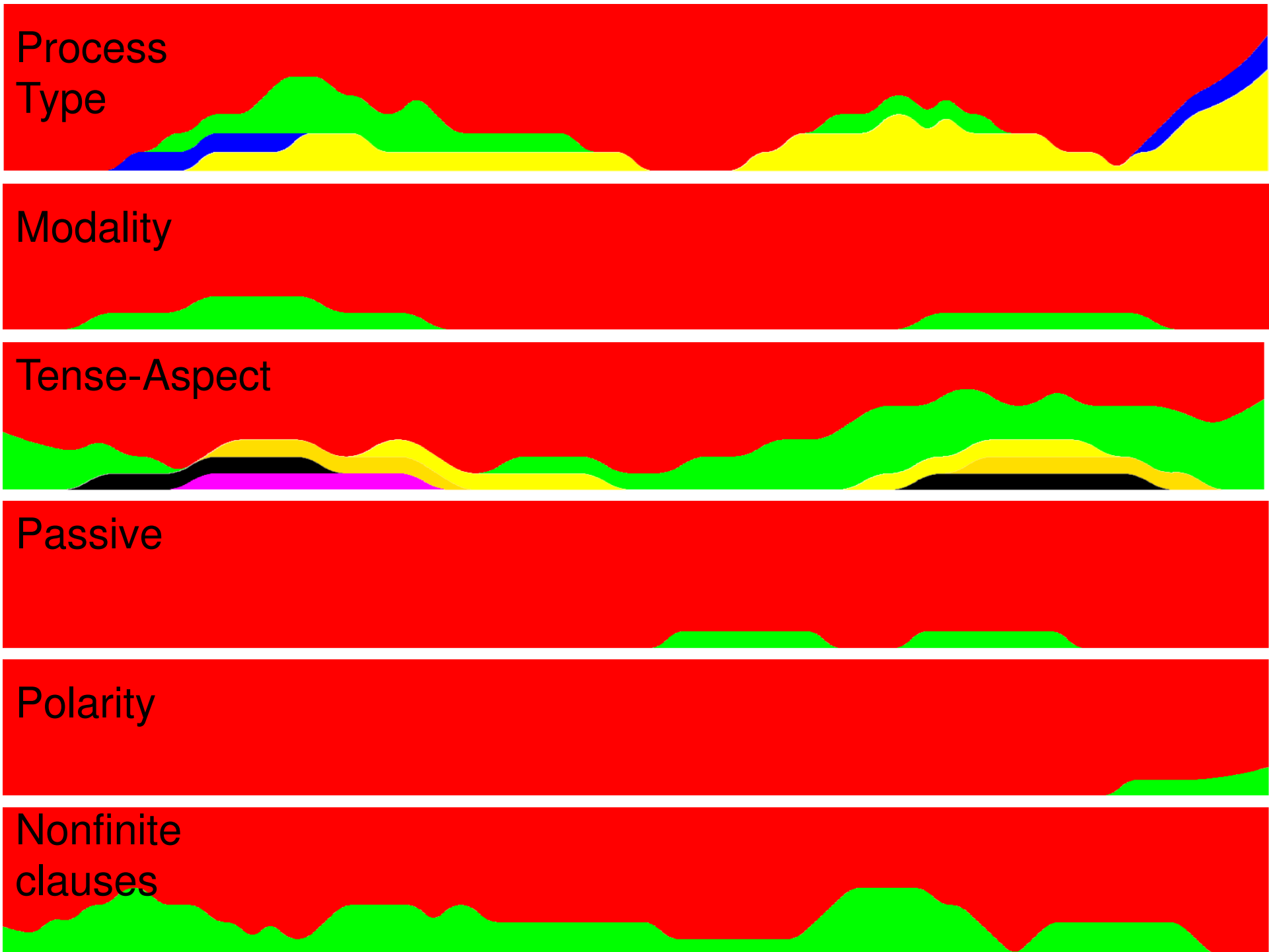
It was the night before the day fixed for his coronation, and the young King was sitting alone in his beautiful chamber. His courtiers had all taken their leave of him, bowing their heads to the ground, according to the ceremonious usage of the day, and had retired to the Great Hall of the Palace, to receive a few last lessons from the Professor of Etiquette. Some of them who had still quite natural manners, which in a courtier is, I need hardly say, a very grave offence.

The lad - for he was only a lad, being but sixteen years of age - was not sorry at their departure. He had flung himself back with a deep sigh of relief on the soft cushions of his embroidered couch, lying there, wild-eyed and open-mouthed, like a brown woodland Faun.

And, indeed, it was the hunters who had found him, coming upon him almost by chance as, bare-limbed and pipe in hand, he was following the flock of the poor goatherd who had brought him up, and whose son he had always fancied himself to be. The child of the old King's only daughter by a secret marriage with one much beneath her in station - a stranger, some said, who, by the wonderful magic of his lute-playing, had made the young Princess love him; while others spoke of an artist from Rimini, to whom the Princess had shown much, perhaps too much honour, and who had suddenly disappeared from the city, leaving his work in the Cathedral

Process
Type

Modality

Tense-Aspect

Passive

Polarity

Nonfinite
clauses

# 5. Conclusions

- Aim was to explore alternative means of visualising patterns in individual texts

- Firstly, split the problem into two issues:

  - What to visualise? What aspect of the text does one wish to explore?

  - How to Visualise? How best to visualise the data to bring out any pattern in the aspect.

- Two forms of data looked at:

  Lexis (keyness, frequency, subjectivity)

  Syntax (keyness, frequency)

- Several visualisations looked at:

  - Table, wordclouds, color-coded text, heat diagrams, textstreams

# 5. Conclusions

- Each visualisation of an analysis has advantages and disadvantages

    – e.g., showing importance vs. showing change through a text

- Good corpus software should provide a range of tools for visualising the data it can extract from the text.

- It is then up to the corpus linguist to choose the appropriate view for their ends.

# 5. Conclusions

- Patterns of a text are not just in lexis

- Recent advances in syntactic parsing allow us to explore the syntactic patterns in our own corpora without doing manual annotation.

- Some of the visualisations used for words can be used also for syntax.

- I showed Tag clouds, which syntactic features are used more in a particular text than in a reference corpus.

- Syntactic choices through a text can also be seen by color coding the text for specific tags.

- "heat" diagrams can show where in a text a particular feature is used.

- TextStreams show where syntactic   alternatives tend to appear in a given text.

# 5. Conclusions

- All of these visualisations have been implemented in UAM CorpusTool (syntax, English only).

- Use can explore different aspects of the text using each of the visualisations (where it makes sense) and choose the one that best reveals the pattern of interest.

- Other aspects of text, and other visualisations, to be added.