

# Macro-trends in research on the central dogma of molecular biology

Sepehr Ehsani

Department of Laboratory Medicine and Pathobiology, and Tanz Centre for Research in Neurodegenerative Diseases, University of Toronto, Toronto, Ontario M5S 3H2, Canada

*Present address:* Whitehead Institute for Biomedical Research, and MIT Computer Science and Artificial Intelligence Laboratory, Cambridge, Massachusetts 02142, United States

ehsani@csail.mit.edu

10 January 2013; revised 7 October 2013

## ABSTRACT

*The central dogma of molecular biology, formulated more than five decades ago, compartmentalized information exchange in the cell into the DNA, RNA and protein domains. This formalization has served as an implicit thematic distinguisher for cell biological research ever since. However, a clear account of the distribution of research across this formalization over time does not exist. Abstracts of >3.5 million publications focusing on the cell from 1975 to 2011 were analyzed for the frequency of 100 single-word DNA-, RNA- and protein-centric search terms and amalgamated to produce domain- and subdomain-specific trends. A preponderance of protein- over DNA- and in turn over RNA-centric terms as a percentage of the total word count is evident until the early 1990s, at which point the trends for protein and DNA begin to coalesce while RNA percentages remain relatively unchanged. This term-based census provides a yearly snapshot of the distribution of research interests across the three domains of the central dogma of molecular biology. A frequency chart of the most dominantly-studied elements of the periodic table is provided as an addendum.*

## MAIN TEXT

The central dogma of molecular biology, consisting of the tripartite alignment of DNA, RNA and protein, was formalized in 1970 [1] (although first put forth in 1958) and has encapsulated the main domains of focus in biological research. Although the boundaries between these domains have become ever more intangible with discoveries such as DNA-protein complexes in the form of nucleosomes [2] or RNA-protein entities such as the ribosome [3], a given research project in the molecular biology laboratory can still be categorized as mainly emphasizing or at least partially involving one of the components of the central dogma more than the remaining two constituents.

It would be informative to determine the extent to which the field of molecular biology has emphasized DNA, RNA and protein as its focus of attention and how the emphasis might have changed over time in the form of a census. Text-mining could be a suitable approach to analyze such patterns in the published literature [4,5]. To that end, 3,767,522 abstracts of all articles containing the word “cell” or “cells” in any field or as a Medical Subject Heading (MeSH) were extracted from the Medical Literature Analysis and Retrieval System Online (MEDLINE) database for a 37-year period from 1975 to 2011, inclusive. Publications indexed prior to 1975 do not necessarily reflect the expected trajectory of increasing annual publication numbers and were not included for the purpose of this study. An in-house word-frequency Java program was used to determine the annual occurrence of a total of 100 DNA-, RNA- and protein-centric terms and their plural/adjectival derivatives, in addition to the total yearly word counts (**Table S1**). Although some of the selected search terms will inevitably capture minor aspects of a domain other than the one assigned, such crossovers will be negligible in the scale and theme of the current search. Moreover, to ensure consistency across different abstracts, the search terms were single words only, and referred to biological entities (e.g., genome) rather than techniques, processes or phenomena (e.g., sequencing, transcription, etc.).

The total number of words, after adjustments for the abstract number and PubMed ID line, was 939,309,404, equalling approximately 250 words per abstract across the 37-year period (**Fig. 1**). The

slope of the total word count, however, is greater than that of the total publications per year, indicating a general trend towards longer abstracts. Whereas the average abstract length from 1975 to 1979 is 180, it increases to 271 words between 2007 and 2011.

The cumulative percentage of protein-centric terms relative to the total annual word count in abstracts from 1975 is 0.52%, whereas it is 0.32% for DNA-centric and 0.13% for RNA-centric terms (**Fig. 2A**). Protein and DNA terms show an increase until the early 1990s, followed by a decline in the late 1990s for protein and mid-2000s for DNA. The trend for RNA shows a relatively slighter variation in the 37-year period. In 2011, the relative percentages for protein, DNA and RNA had changed to 0.47%, 0.46% and 0.16%, respectively. It should be noted that although the number of publications in each domain has increased substantially every year, the percentage of search terms per total words indicates the relative focus of the field of molecular biology on each part of the central dogma. Overall, it is evident that the protein component of the central dogma of molecular biology received greater attention in research in the 1970s, followed by DNA and RNA. With a gradual decline in protein-centric percentages and relatively constant increase in DNA terms, the protein and DNA domains appear to have garnered equal focus in the 2000s. The relative attention on the RNA domain seems to have remained less fluctuating.

It will be useful to determine if the trend observed for each domain is mainly carried by one or more subdomains. 'Gene' and related terms are the predominant trendsetters for DNA-centric words (**Fig. 2B**). Other than 'genome' and 'epigenetics' (and related terms), which show a slow rise from early-2000s onward, all other terms have a negative or near-zero slope. The trend for RNA-centric words seems to be mainly carried in phases by three different subcategories: from 1975 to the mid-1980s by 'RNA', from the mid-1980s to the early 2000s by 'mRNA' (messenger RNA), and from then on by 'siRNA' (small interfering RNA) and 'miRNA' (microRNA) (**Fig. 2C**). The main term in the protein-centric trend is, unsurprisingly, 'protein' and related terms (**Fig. 2D**). However, the key contributors to the shape of the trend in the 1970s and 1980s are 'enzyme' and 'antibody', respectively, both of which demonstrate a steady decline in subsequent decades. The 'proteome' subdomain shows a relatively small rise in the past decade.

All 100 search terms account for an average of 1.15% of the total words in the abstracts with a standard deviation of 0.11%. If the three constituents of the central dogma encompass the principal avenues of information flow in the cell, it might then be plausible to assume that the cumulative percentage of all DNA-, RNA- and protein-related terms should be constant on a year-by-year basis, even in cases of new discoveries within the dogma which spawn new research directions and publications. Although the current census appears to corroborate this assumption to an extent, the sums of all percentages span a minimum of 0.95% in 1977 and a maximum of 1.32% in 1995 (**Fig. S1**). It is possible to speculate that publications analyzed in this study have at different times relied on biological phenomena outside of the tripartite axis (e.g., phospholipids) to complement their explanation of a given aspect of the cell's information system.

To control for the accuracy of the search algorithm, the frequency of the articles 'a', 'an' and 'the' were calculated (**Fig. S2**). The indefinite articles show a consistent percentage throughout the search period, with 'a' and 'an' representing 1.5% and 0.3% of the total words, respectively. Interestingly, the usage of the definite article 'the' shows a negative linear trend, with a decrease from 6.0% in 1975 to 3.8% in 2011.

Finally, to determine the relative focus of the field of molecular and cellular biology on the natural elements, the frequencies of occurrence of the 118 elements in the periodic table were calculated (**Fig. S3**). Calcium, oxygen, lead, iron, sodium, carbon, zinc, hydrogen, gold, copper and potassium represented the most frequently mentioned elements at present, respectively (**Fig. S3A**). Relative to other elements, calcium and sodium demonstrate a decrease in occurrence from the late 1980s to 2011 (**Fig. S3B**).

## REFERENCES

1. Crick F (1970) Central dogma of molecular biology. *Nature* 227: 561-563.
2. Kornberg RD (1974) Chromatin structure: a repeating unit of histones and DNA. *Science* 184: 868-871.
3. Hoagland MB, Stephenson ML, Scott JF, Hecht LI, Zamecnik PC (1958) A soluble ribonucleic acid intermediate in protein synthesis. *J Biol Chem* 231: 241-257.
4. Andrade MA, Bork P (2000) Automated extraction of information in molecular biology. *FEBS Lett* 476: 12-17.
5. Jensen LJ, Saric J, Bork P (2006) Literature mining for the biologist: from information retrieval to biological discovery. *Nat Rev Genet* 7: 119-129.

## FIGURE LEGENDS

**Figure 1. Overview of analyzed abstracts.** The total number of publications (bar graph, left vertical axis) shows a steady rise from tens of thousands in the late 1970s to hundreds of thousands in the late 2000s. The total number of words for each year (line graph, right vertical axis), however, increases more rapidly over the same time period and indicates a trend toward longer abstract lengths.

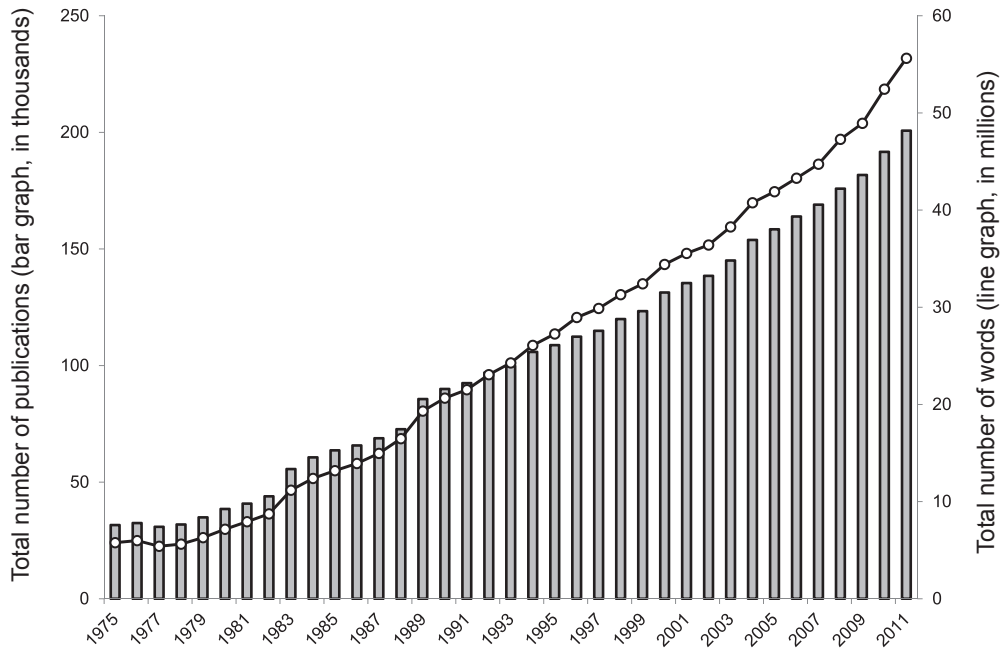
**Figure 2. Central dogma domain trends from 1975 to 2011.** Protein-centric research appears to have dominated the field of molecular biology in the late 1970s, followed by DNA-centric and RNA-centric research (**A**). With a steady rise in DNA-centric terms and decline in protein-centric terms, the DNA and protein domains converge on equal percentages relative to total words in the 2000s. The breakdown for individual trends demonstrating the contribution of each subdomain appears in **B**, **C** and **D**.

**Figure S1. Annual sums of percentages.** A revised **Fig. 2A** includes a cumulative line graph of the yearly DNA-, RNA- and protein-centric percentages, which range from 0.95% (1977) to 1.32% (1995).

**Figure S2. Control graph of the frequency of definite/indefinite articles.** The indefinite articles 'a' and 'an' demonstrate a steady frequency over the period 1975-2011, whereas the definite article 'the' decreases from 6.0% to 3.8% over the same period.

**Figure S3. Frequency of the elements.** The absolute (**A**) and relative (**B**) occurrences of elements 1-118 were calculated. Two alternative spellings were considered for four elements: aluminium/aluminum, phosphorus/phosphorous, sulfur/sulphur and caesium/cesium. The elements francium, berkelium, einsteinium, mendelevium, nobelium, lawrencium, rutherfordium, dubnium, seaborgium, bohrium, hassium, meitnerium, darmstadtium, roentgenium, copernicium, ununtrium, flerovium, ununpentium, livermorium, ununseptium and ununoctium did not occur in any abstract.

Figure 1



**Figure 2**

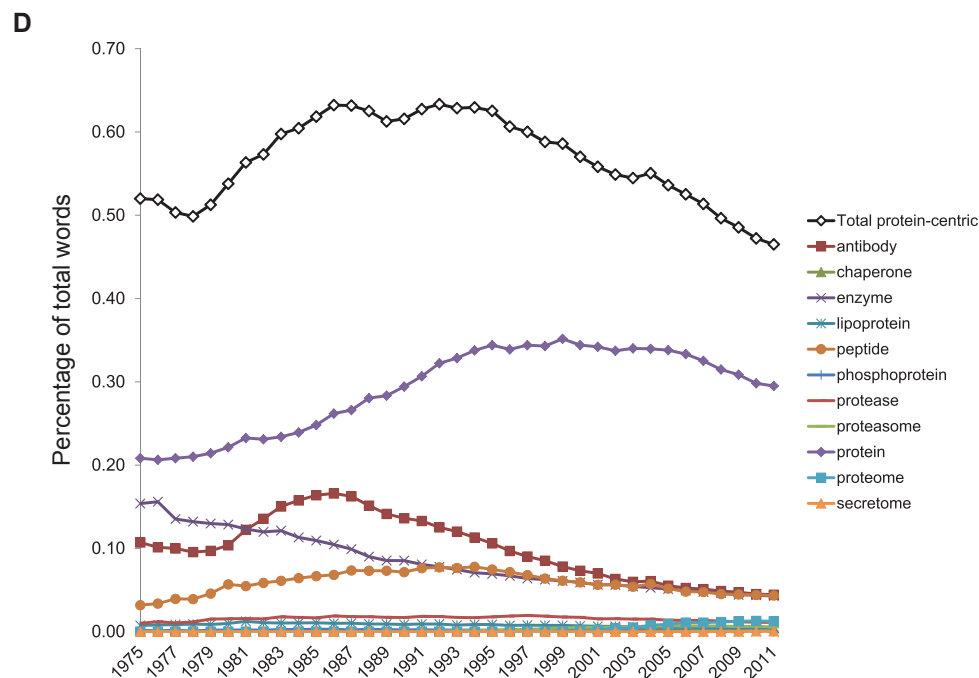
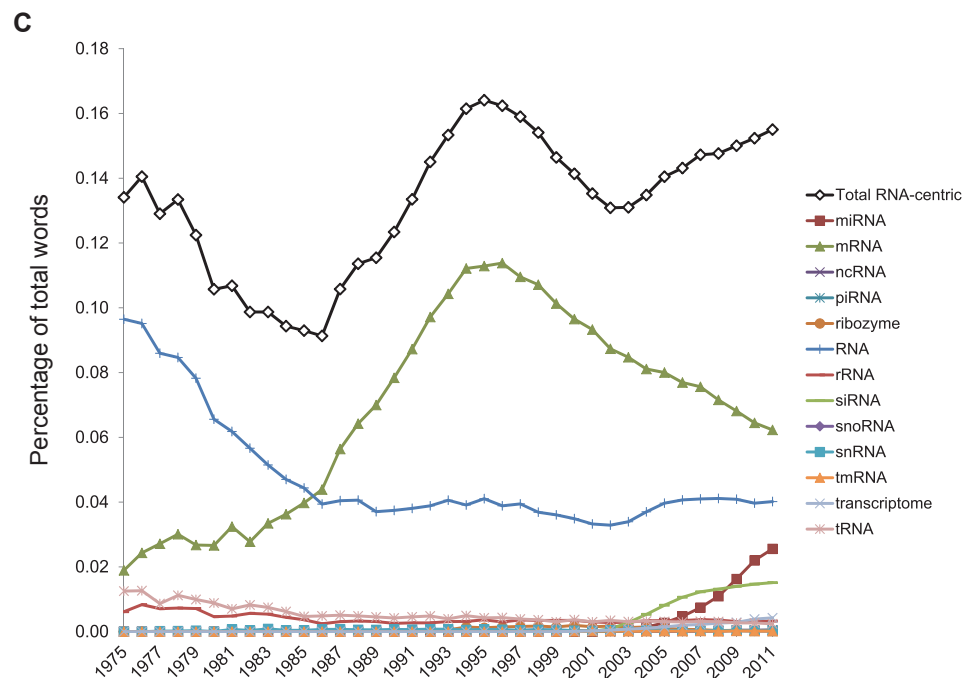
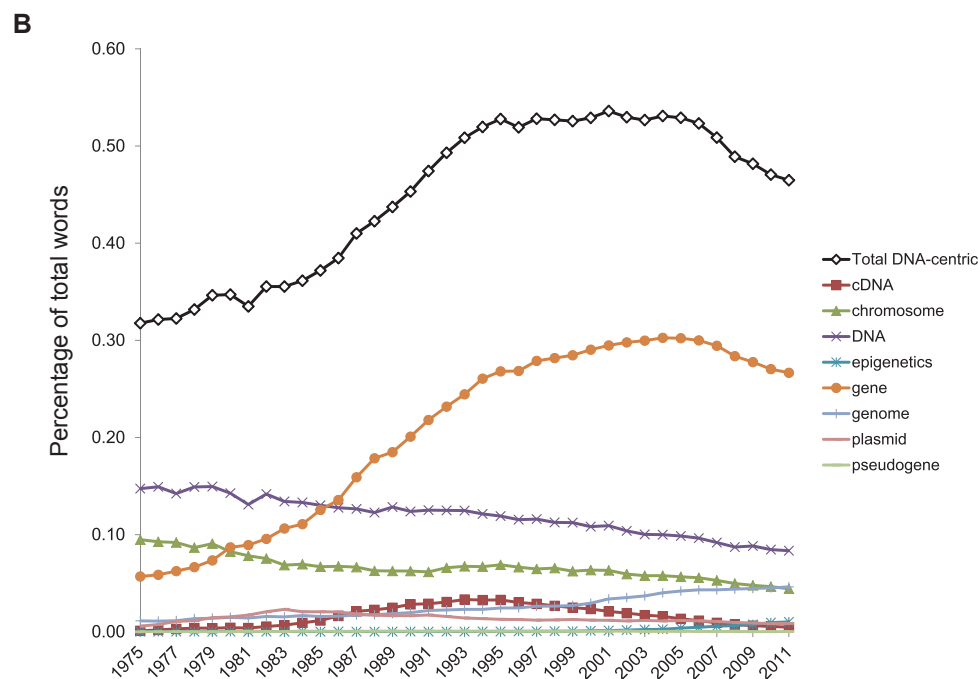
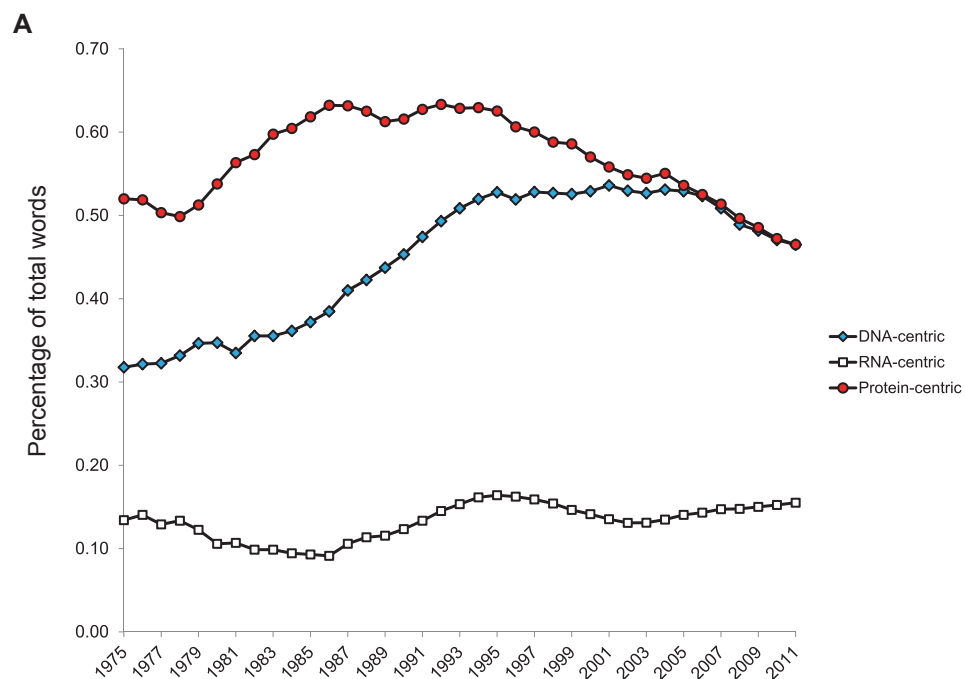




Figure S1

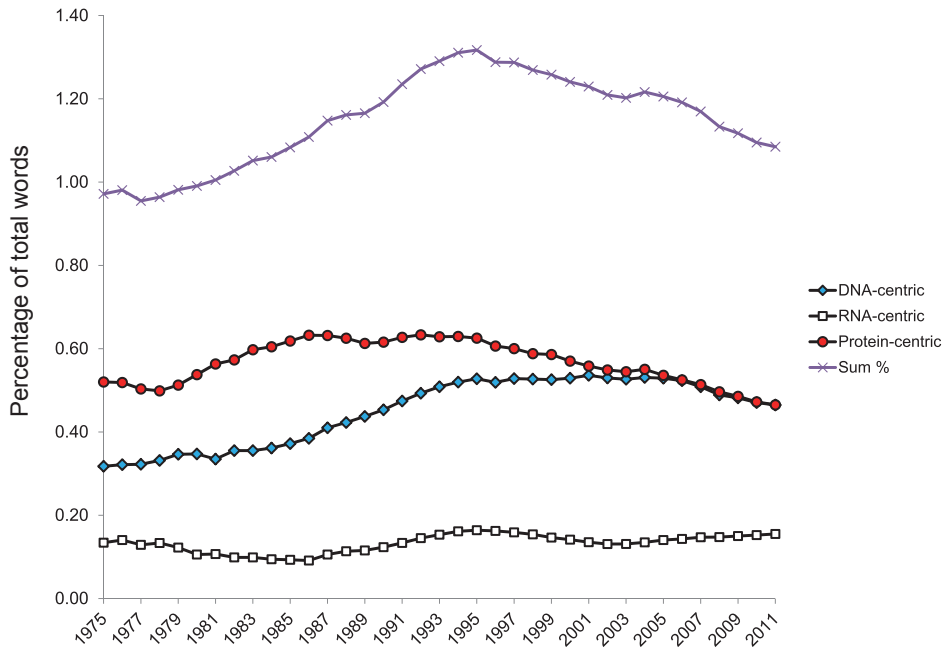


Figure S2

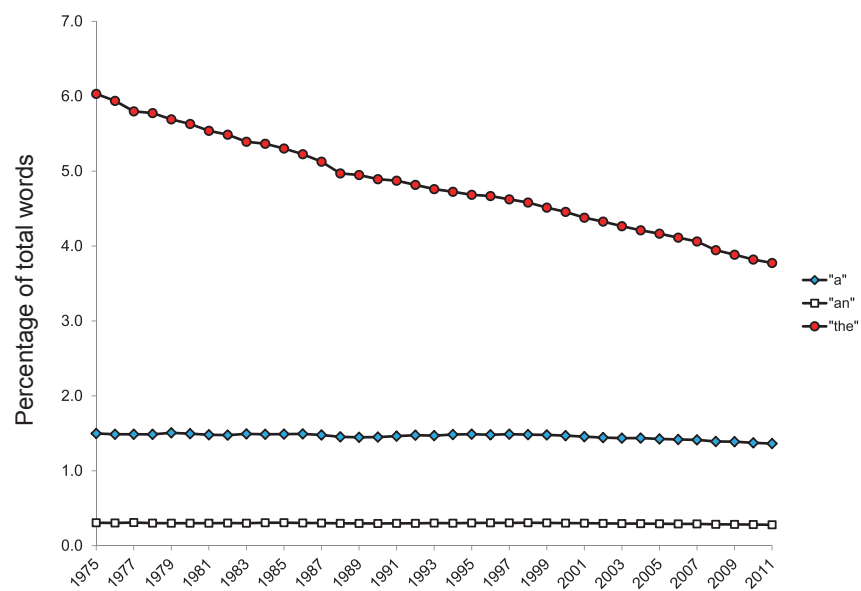
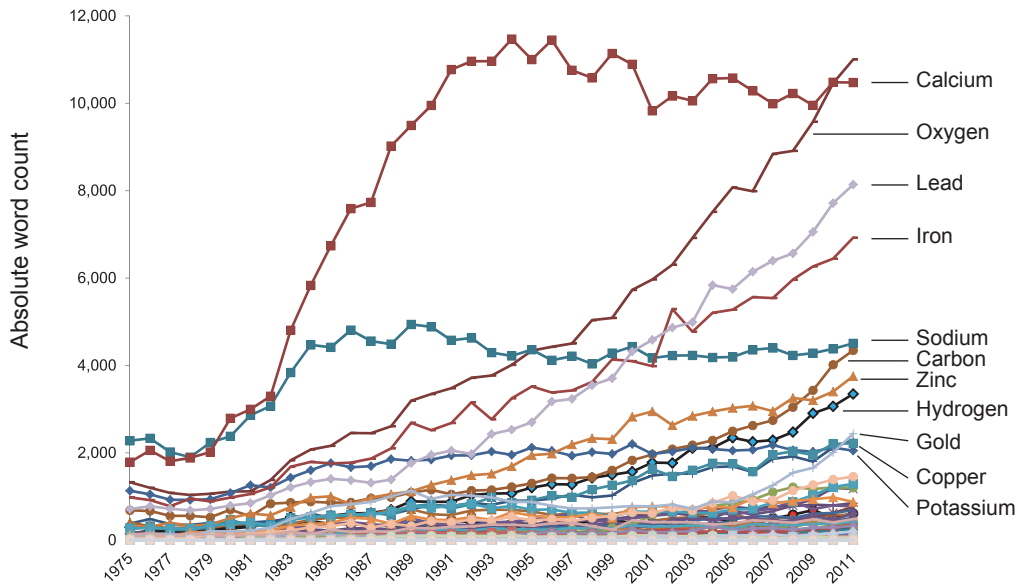




Figure S3

A



B

