

Available online at www.sciencedirect.com

ScienceDirect

Journal homepage: www.elsevier.com/locate/cortex

Special issue: Editorial

Language, computers and cognitive neuroscience

Peter Garrard^{a,*} and Brita Elvevåg^{b,c,**}^a Neuroscience Research Centre, Institute of Cardiovascular and Cell Sciences, St George's, University of London, Cranmer Terrace, London, UK^b Psychiatry Research Group, Department of Clinical Medicine, University of Tromsø, Norway^c Norwegian Centre for Integrated Care and Telemedicine (NST), University Hospital of North Norway, Tromsø, Norway

It is not tissue microstructure or functional capability that sets the human brain apart from other organs and systems, but its organisational complexity, and to understand the brain at this level remains one of the great scientific challenges of our age. There is no doubt that computation will prove central to the endeavour, both as a framework for understanding, and a medium for simulating, cognition and its myriad disorders. The power and interconnectedness of modern computing hardware are now being exploited in some of the largest and most ambitious studies of cognition ever undertaken [*'Head Start'* (Editorial comment) *Nature*, 2013]. The availability of supercomputing power also opens up the related possibility of exploiting novel information sources that are too large and complex to be captured, organised or analysed using conventional approaches – a resource that, over recent years, has come to be known as 'big data'. The McKinsey Global Institute's 2011 report on this phenomenon is entitled *Big data: The next frontier for innovation, competition, and productivity* (Manyika et al., 2011). The authors show how big data generate value in healthcare, public services, retail and manufacturing. Among our ambitions for this *Cortex* special issue is that it will help to make the case for cognitive neuroscience to be added to the list.

How might big data contribute to the goals of understanding healthy and disordered brains in ways that span Marr's three 'levels of analysis'? (Marr, 1982) (See also Poggio's recent update on this framework, which is available in full at: <http://cbcl.mit.edu/publications/ps/MIT-CSAIL-TR-2012-014.pdf>.) In a world dominated by digital technology,

people contribute to the store of big data simply by going about their daily lives. The focus of interest in the resulting, and constantly growing, body of information will naturally vary: for the business community the behaviour, choices and preferences of users and customers will be critical to the goal of maximising profits, while government and the public sector must aim to formulate indices of economic value and social outcome in order to maximise the efficient use of limited resources. Meanwhile, science has both benefited from and pioneered the understanding of huge datasets and data streams, including those related to particle physics, genomics and climate science – fields that generate quantities of data measured in petabytes ($\times 10^{15}$ bytes) per year (Doctorow, 2008).

It is inevitable that the information people generate as they go about their daily lives will hold some value for cognitive neuroscientists, particularly those who emphasise the importance of 'ecological validity' in the interpretation of behavioural data (Cohen, 1996; Neisser, 1991). We have no interest in reopening any of the wounds inflicted (by both sides) in the debate on the relative merits of everyday memory and traditional laboratory research. Yet few people with a scientific interest in learning and memory would dismiss out of hand a detailed and cumulative record of (for example) all the movements, interactions and web searches carried out by large populations of individuals over a number of years. Although the level of intrusion that would be required to generate such a dataset on private citizens is hardly desirable, there is less cause for squeamishness when one considers the

* Corresponding author. Neuroscience Research Centre, Institute of Cardiovascular and Cell Sciences, St George's, University of London, Cranmer Terrace, London SW17 0RE, UK.

** Corresponding author. Psychiatry Research Group, Department of Clinical Medicine, University Hospital of North Norway – Åsgård, Postbox 6124, 9291 Tromsø, Norway.

E-mail addresses: pgarrard@sgul.ac.uk, peter.garrard@gmail.com (P. Garrard), brita@elvevaag.net (B. Elvevåg).

<http://dx.doi.org/10.1016/j.cortex.2014.02.012>

0010-9452/© 2014 Elsevier Ltd. All rights reserved.

benefits that could accrue to closed communities, be they real (e.g., work environments or care homes) or virtual (e.g., patient groups with internet connections and/or access to clinical care via a telemedicine programme). Yet there is much practical and ethical ground to move before such data become relevant and usable.

Fewer limitations apply to language data in the form of naturally produced samples of spoken or written language: collection and recording have been taking place for hundreds of years in the form of handwritten and printed texts, audio recording, and most recently the hundreds of millions of digital communications (blogs, tweets, emails and text messages) that are produced each day by an ever more digitally interconnected public. Many of these sources are the product of undirected and spontaneous cognitive activity in single individuals, often with the intention of public communication. In addition, there is a sizeable body of clinical data representing the output of more focused neurocognitive activity in various clinically defined groups (of which, more later). All can be considered in the light of a multitude of dimensions, some of them simple, others reflecting more complex attributes of the symbolic systems in which they are represented. The widespread availability of fast, high capacity, desktop computers means that large volumes can be represented and stored in a digital text format.

Nonetheless, the problem of how to make large datasets tractable, to organise and use them in informative ways remains common to all the enterprises – scientific, technical and commercial – that we have considered so far. Previous attempts to extract meaning from huge datasets have relied on a diverse range of ‘data mining’ techniques, including dimension reduction, information theory, and statistical machine learning – approaches that are represented in a number of the contributions to this special issue. Even simple approaches such as proportional word-counts, however, can produce strikingly informative results, particularly when applied to very large datasets. A leading source of both data and analytical tools is Google: Google Books contains digitally encoded texts of a large (and ever increasing) proportion of all the books ever published; the Google n-gram viewer <https://books.google.com/ngrams> will plot the change in proportional frequency of any word (unigram) or phrase (n-gram) in books published between the years 1800 and 2000. In a series of fascinating explorations of the data, Michel et al. (2011) reported a selection of instances in which social and cultural evolution and major historical events were reflected in lexical frequency trends. The approach offers limitless possibilities for further exploration, and it is to be hoped that the interdisciplinary nature of cognitive neuroscience will prompt experts from disciplines such as statistics and computer science to modify and add to the analytical armamentarium.

Even if the ‘what?’ and ‘how?’ of large scale language analysis could be fully addressed, we would still be left with the question that even the most rarefied scientific disciplines must nowadays address: ‘to what end?’ We contend that the contributions to this special issue provide a wealth of justifications, predominantly clinical, but also theoretical. Among the latter are the contributions of Montemurro (2014) and Voorspoels et al. (2014). The former advances the idea that the inherent order detectable in the long range co-occurrence of

words in texts written in different languages (relative entropy) should be considered a candidate for a quantitative linguistic universal – a bold and testable hypothesis. The latter explores the pitfalls and limitations of the clustering method in arguing for distorted semantic structure in cognitive neuropsychology. Valle-Lisboa, Pomi, Cabana, Elvevåg, and Mizraji (2014) adopt a neurocomputational modelling approach to explore the links between matrix associative memory models and language processing and production, creating a system for exploring how disruptions in connectivity between the underlying representations of concepts can result in various forms of disorganized speech.

Clinically based studies draw on a wide-ranging series of data associated with language change over the course of normal ageing (Ferguson et al., 2014) and tenure of political office (Garrard, Rentoumi, Lambert, & Owen, 2014), as well as linguistic features of cerebral functional disorders including Alzheimer’s disease, primary progressive aphasia (Garrard, Rentoumi, Gesierich, Miller, & Gorno-Tempini, 2014), and schizophrenia. These studies are made possible by the fact that communication is a high-level neurocognitive function providing a rich and extemporaneous dataset that reflects the state of numerous interacting neural and cognitive processes. If assayed appropriately, therefore, communication affords a unique and sensitive window into a person’s state of mental and cognitive health.

As authors, we welcome exposure of our research to the more than usually diverse readership that the interdisciplinary theme of this special issue will attract. As editors, we were struck by the multiplicity of ways in which computer-assisted analysis of large language datasets could contribute to the understanding of brain disorders. Pakhomov and Hemmy (2014) took a large database of verbal fluency responses collected as part of the Wisconsin Nun Study (Snowdon et al., 1996) and interrogated the data for response clusters and switching behaviours using an automated measure of relatedness derived from latent semantic analysis (LSA). Originally conceived as a statistical approach to the acquisition and representation of meaning (Landauer & Dumais, 1997), LSA uses a vector space representation of the words and contexts occurring in large numbers of digitised texts, such that the distance between vectors can be used as a metric of the semantic similarity between the words and/or contexts. This property allows a number of robust measurements to be made in novel text or discourse samples, including those obtained from different patient groups.

Hoffman, Meteyard, and Patterson (2014) use the neighbourhood density of items in a semantic space to derive a measure of ‘semantic diversity’ characterising the vocabulary of patients with conceptual degradation (semantic dementia). Several studies use LSA to examine the properties of discourse samples obtained from patients with schizophrenia. Two papers (those by Holshausen, Harvey, Elvevåg, Foltz, & Bowie, 2014; Tagamets, Cortes, Griego, & Elvevåg, 2014) report correlations between LSA derived measures of patient discourse and other validated functional measures, namely clinical and psychometric indices, and task-related fMRI patterns. A third (Rosenstein, Diaz-Asper, Foltz, & Elvevåg, 2014) examines the effect of a latent semantic variable and a syntactic characteristic to examine the effects of these features on prose recall

Download English Version:

<https://daneshyari.com/en/article/7315495>

Download Persian Version:

<https://daneshyari.com/article/7315495>

[Daneshyari.com](https://daneshyari.com)