

<https://helda.helsinki.fi>

Semiotically-grounded distant viewing of diagrams : insights from two multimodal corpora

Hiippala, Tuomo

2022-05-25

Hiippala , T & Bateman , J A 2022 , ' Semiotically-grounded distant viewing of diagrams : insights from two multimodal corpora ' , Digital scholarship in the humanities , vol. 37 , no. 2 , pp. 405-425 . <https://doi.org/10.1093/lc/fqab063>

<http://hdl.handle.net/10138/344232>

<https://doi.org/10.1093/lc/fqab063>

cc_by

publishedVersion

Downloaded from Helda, University of Helsinki institutional repository.

This is an electronic reprint of the original article.

This reprint may differ from the original in pagination and typographic detail.

Please cite the original version.

Semiotically-grounded distant viewing of diagrams: insights from two multimodal corpora

Tuomo Hiippala 

Department of Languages, University of Helsinki, Helsinki, Finland

John A. Bateman

Department of English, Universität Bremen, Germany and
Department of Linguistics, Universität Bremen, Germany

Abstract

In this article, we argue for the benefits of combining large-scale analyses of visual materials currently pursued within digital humanities with insights from multimodality research, which is an emerging discipline that studies how human communication relies on appropriate combinations of expressive resources. We show that concepts developed within the field of multimodality research provide appropriate metadata schemes for various modes of expression in large corpora and datasets. We illustrate the proposed approach using a common mode of expression, diagrams, and analyse two recent multimodal diagram corpora using statistical and computational methods. Our results suggest that multimodally-motivated metadata schemes can provide a robust foundation for computational analyses of large corpora and datasets. Even if a corpus or dataset is not designed to support full-blown analyses of multimodal communication, our results imply that multimodality theory can still be used to impose tighter analytical control over a variety of visual materials.

Correspondence:

Tuomo Hiippala, Department of Languages, University of Helsinki, Yliopistonkatu 4, 00100 Helsinki, Finland.

E-mail:

tuomo.hiippala@helsinki.fi

1 Introduction

Whether taking place via an external medium or in face-to-face interaction, communication is naturally multimodal: that is, making and exchanging meanings involve combining multiple modes of expression in a coordinated, goal-oriented manner. There is currently growing interest in multimodal communication across various fields of research, including the digital humanities. It is then natural to consider more closely whether contemporary theories of multimodality can support the kinds of large-scale analyses commonly

pursued in digital humanities and if so, to what extent. Although the field of multimodality is increasingly oriented towards empirical analysis, compiling multimodal corpora to support such analyses is still highly labour-intensive. More extensive use of computational techniques is thus a clear priority.

In this article, we consider the potential benefits of combining contemporary accounts of multimodality, computational methods, and research orientations from digital humanities with respect to one extremely common mode of expression, namely diagrams. Diagrams are found everywhere, and their structure

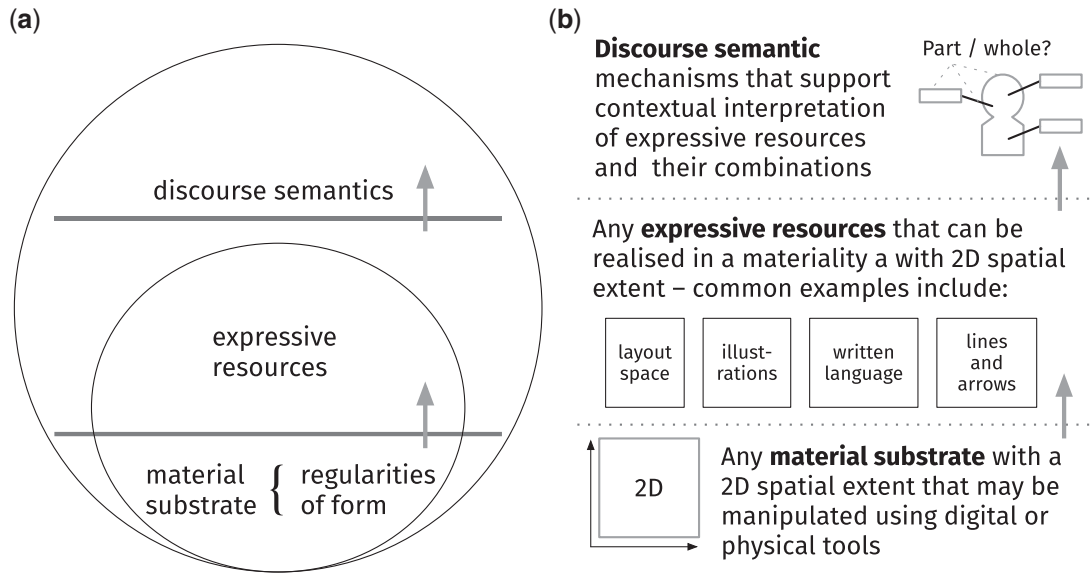


Fig. 1 The concept of a semiotic mode and its application to diagrams. (a) A theoretical model of a semiotic mode. (b) A characterization of the diagrammatic mode.

language, illustrations, and photographs may *naturally* co-occur with each other in diagrams, thus avoiding committing to arbitrary divisions between ‘verbal’ and ‘visual’ or ‘text’ and ‘image’ (Bateman, 2014). This perspective is obviously carried over to mass media, which regularly deploy multiple semiotic modes (Bateman *et al.*, 2017, p. 124). Finally, discourse semantics guides the interpretation of expressive resources and their combinations in context. For diagrams, resolving the resulting discourse relations relies on formal cues such as spatial placement of elements or connections realized using lines and arrows in combination with world knowledge (Watanabe and Nagao, 1998; Alikhani and Stone, 2018).

This brief description illustrates the extent to which modern multimodality theory can explicate how semiotic modes operate quite generally. Describing the characteristics of each stratum of a given semiotic mode—that is, material substrate, expressive resources, and discourse semantics—is then an issue demanding empirical research. Conversely, this also shows just how much complexity is missed when operating with pretheoretical distinctions such as ‘text’ and ‘image’. Although computational analyses

of page-based media are already advancing beyond such dichotomies, as Wevers and Smits (2020) have shown by training convolutional neural networks to distinguish between instances of illustrations, photographs, and other semiotic modes in historical newspapers, we argue that advancing this effort within digital humanities will benefit still further from the input of multimodality theory.

As a form of ‘applied semiotics’ that seeks a close relationship between theory and data (Bateman and Hiippala, 2021), multimodality theory is well-positioned to provide a foundation for characterizing the diverse range of communicative artefacts and situations studied within digital humanities (Bateman, 2017; Hiippala, 2021). Theories of multimodality have already been used to guide the application of computational methods to both filmic (Bateman *et al.*, 2016) and page-based media (O’Halloran *et al.*, 2018), but much remains to be done in terms of applying computational methods in a way that respects the complexity of multimodal communication. While we are not suggesting that all studies that use computational methods should perform full-blown multimodal analyses for each semiotic mode encountered, we do

encourage the broader application of multimodality theory to determine the analytical granularity appropriate for answering specific research questions. This allows targets of descriptions and their respective granularities to be derived systematically on the basis of developing bodies of theory. There are then both theoretical and practical reasons for adopting this approach when collecting multimodal data at scale.

As noted by [Arnold and Tilton \(2019, p. i4\)](#) above, it is important when constructing larger collections or corpora for computational analyses that appropriate metadata schemes are defined for organizing that data. In computer science, the definition of ‘modality’ (the term preferred over ‘mode’) is strongly aligned with the senses: our ability to see, hear, touch, and use natural language. Distinctions between sensory modalities are then built into the different research fields of computer vision, audio signal processing, and natural language processing, and these are then the sources for corresponding metadata schemes. The resulting fields are as a consequence often confined to their own ‘problem spaces’, even though these are increasingly converging on multimodality in tasks such as machine translation ([Sulubacak et al., 2020](#)). Nevertheless, applying definitions based on sensory modalities continues. In contrast, within humanities-oriented multimodality theory, restricting approaches to sensory channels is now receiving considerable critique because defining modalities ahead of analysis solely on the basis of perceptual properties makes identification of the actual semiotic contributions being made to meaning construction more difficult. As argued extensively in [Bateman et al. \(2017\)](#), semiotic contributions regularly extend across sensory channels and their demarcations need to be teased out empirically: one cannot assume their individual characteristics in advance ([Bateman, 2011, p. 17–18](#)). It is crucial to open up a two-way communication channel between the semiotic distinctions being made and their supporting material distinctions, rather than assuming that sensory perception alone will result in appropriate segmentations.

These challenges and limitations are made fully evident by our characterization of the diagrammatic mode in [Fig. 1b](#). First, diagrams are clearly not aligned with a single traditional modality as they cross-cut both ‘vision’ and ‘language’. Second, what makes diagrams different from other combinations of a similar

nature, such as photographs with embedded or overlaid text, is left an open question. Although assumptions of similarity concerning expressive resources *within* a single sensory modality are common in computer vision research, where objects of analysis are often reduced to mere carriers of content, this is insufficient. [Haehn et al. \(2019, p. 649\)](#), for example, report that models trained on photographs do not generalize well to diagrammatic representations without further training even though they are both clearly ‘visual’. They consider this finding surprising given prior comparisons between artificial neural networks and the human visual cortex, which assume that *visual perception* suffices for reasoning about both photographs and diagrams.

From the perspective of multimodality theory, however, the differences between photographs and diagrams are rather evident: diagrams differ radically in terms of their expressive resources and discourse semantics. Diagrams are compositional, that is, they can be broken down into component parts, which may be realized using multiple expressive resources and combined into discourse structures that work towards a shared communicative goal. This allows diagrams to represent abstract concepts and phenomena that are not limited to specific slices of time and space, and so stand in strong contrast to photographs (cf. e.g. [Alikhani and Stone, 2018](#); [Greenberg, 2018](#)). This demonstrates how it is always essential to consider material distinctions in terms of the particular semiotic modes they operate with respect to. It is precisely these semiotic modes that deliver appropriate metadata schemes for characterizing corresponding objects of analysis.

3 Insights from Multimodal Diagram Corpora

Having introduced the concept of a semiotic mode and how appropriate metadata schemes may be derived for individual semiotic modes through empirical research, we turn now to examine two recent diagram corpora from this perspective. These corpora originate in two different fields of research, namely artificial intelligence ([Kembhavi et al., 2016](#)) and multimodality research ([Hiippala et al., 2020](#)), but contain

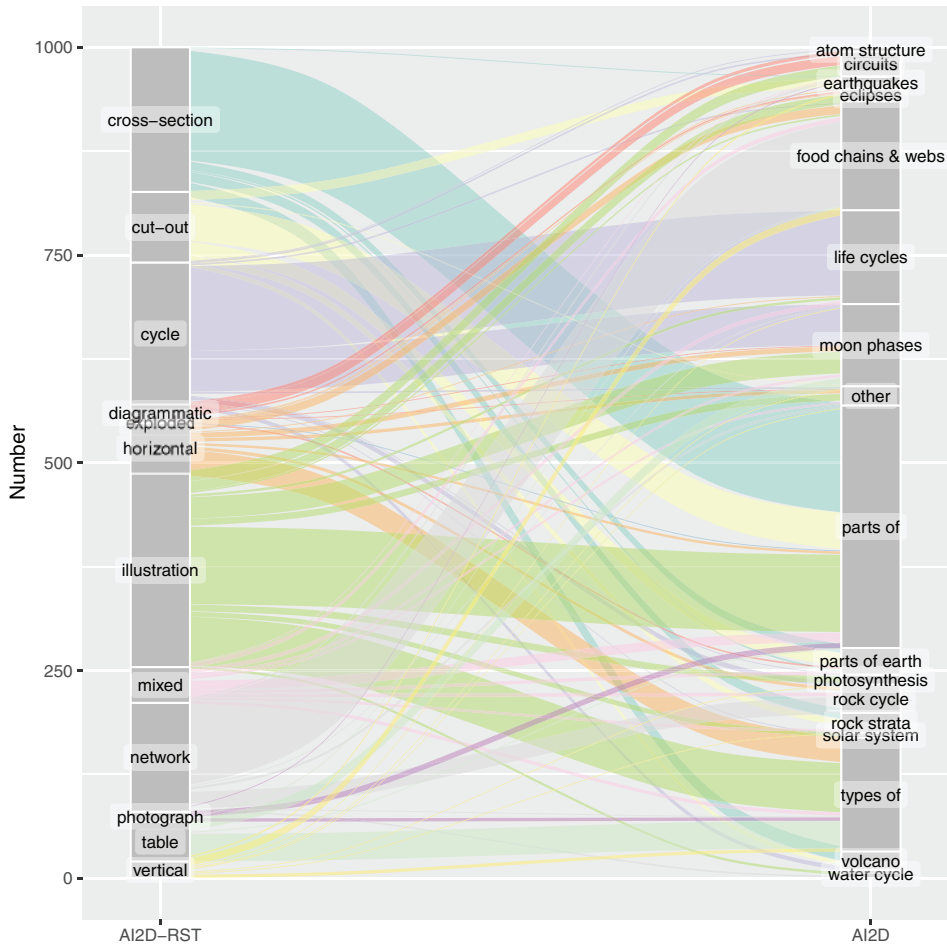


Fig. 6 An alluvial plot mapping the structural categories in AI2D-RST to the semantic categories in AI2D. Y-axis shows the number of diagrams in each category. Diagrams in AI2D-RST that combine multiple categories are labelled as ‘mixed’. Created using *galluvial* 0.12.3 for R 4.0.2.

populating these spatial positions with blobs may evoke an association with cycles, thus operating as a genre cue that encourages the viewer to consider whether such a discourse semantic interpretation holds.

The remaining examples in Fig. 8 show layout patterns for three further structural categories in AI2D-RST concerned with how diagrams represent their depicted objects (Hiippala *et al.*, 2020, p. 8–9). These include *illustration* in Fig. 8b, which covers all forms of depiction at various levels of visual detail from monochrome to colour drawings. *Illustration* is distinguished from *cross-section* in Fig. 8c and *cut-out*

in Fig. 8d based on whether the internal structure of the depicted object is shown by cutting the object in half (*cross-section*) or by removing a part of the object to expose its structure (*cut-out*). As the layout patterns show, *illustrations* are more flexible in their positioning of blobs than *cross-sections* and *cut-outs*, as *illustrations* occasionally depict multiple objects in a single diagram, which causes the blob centroids to spread out. The layout patterns for *cross-sections* and *cut-outs*, in turn, cannot be distinguished from one another based on layout information alone.

This exposes certain limitations of the metadata scheme used to describe diagram elements in the

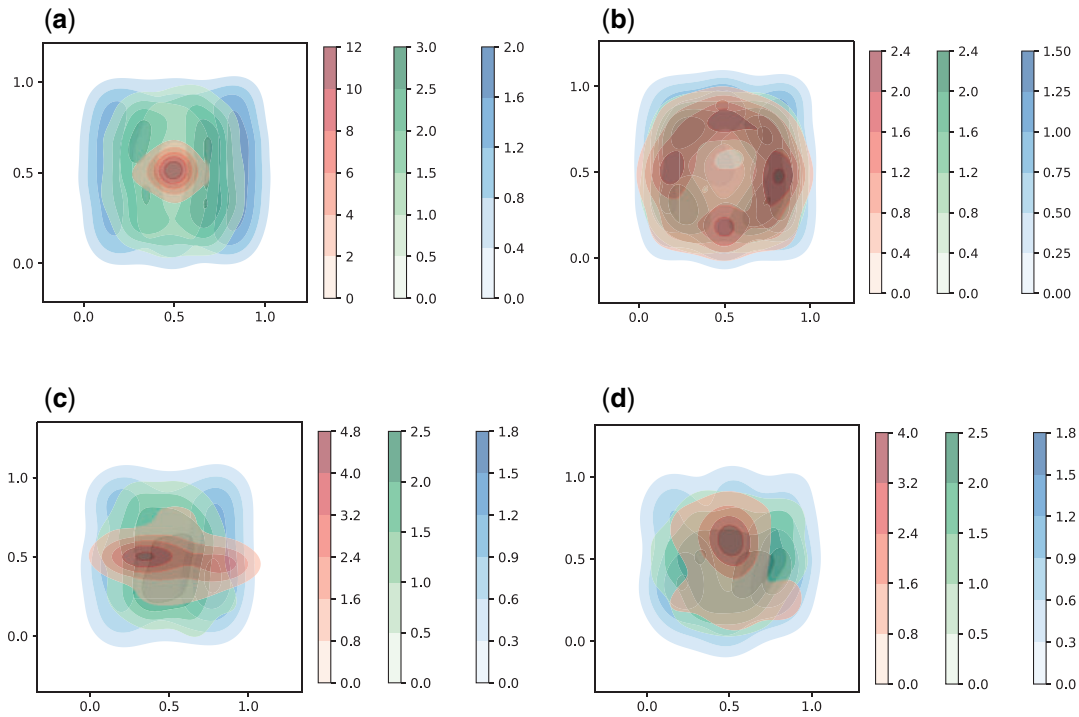


Fig. 7 KDEs for the centroids of text (blue), arrow/line (green), and blob (red) elements for four semantic categories in the AI2D dataset. The coloured bars to the right of the graphs map the colours to the values of probability density function for each diagram element type. Created using *matplotlib* 3.3.0 (Hunter, 2007) and *seaborn* 0.10.1 for Python 3.8.5. (a) Parts of. (b) Life cycles. (c) Rock strata. (d) Volcano.

AI2D dataset, which may be traced back to defining the expressive resources *ahead* of actual analysis (see Section 3.1). Because AI2D classifies all forms of depiction as ‘blobs’, we cannot determine whether categories in AI2D-RST such as *illustrations*, *cycles*, and *cut-outs* prefer to draw on different expressive resources—for example coloured hand-drawn illustrations or monochrome line drawings—for depicting objects and their structure. Put differently, forms of depiction must be described more accurately to identify which expressive resources are being deployed and for which communicative purposes in each diagram category. Capturing these distinctions is crucial because the diagrammatic mode uses such expressive resources to adjust diagrams’ levels of abstraction (Dimopoulos *et al.*, 2003). To exemplify, an animal may be represented by a round circle in a *network* diagram showing its role in a food web, whereas an *illustration* is more

likely to use a lifelike drawing to portray the same animal. Distinguishing between these representations was not possible because in AI2D and AI2D-RST both are classified as ‘blobs’. In the following section, we explore the diversity of expressive resources constituting the category of ‘blobs’ using computer vision methods.

3.4 Unpacking the expressive resources in ‘blobs’

We now offer a more fine-grained description of the expressive resources collectively labelled as ‘blobs’ in the AI2D dataset. To explore which expressive resources are used for depiction in *illustrations*, *cycles*, and *cut-outs*, we extract all diagram elements classified as blobs from the AI2D dataset ($N = 20,937$) and apply the method presented in Fig. 9 to characterize their visual appearance in terms of brightness and texture.

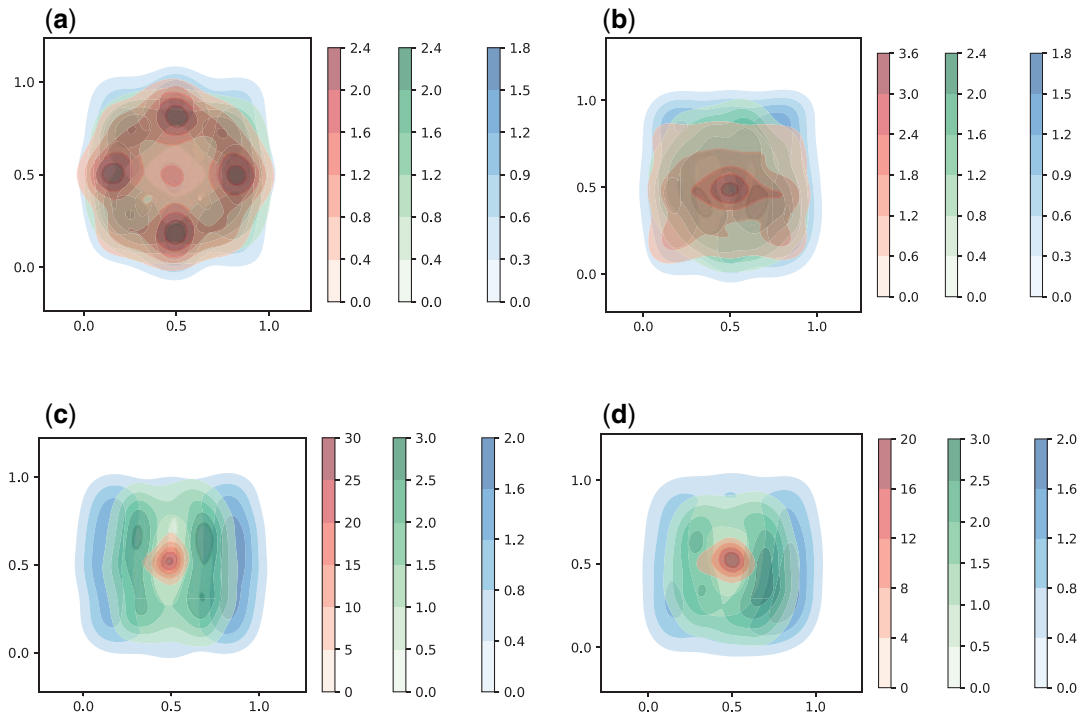


Fig. 8 KDEs for the centroids of text (blue), arrow/line (green), and blob (red) elements for four structural categories in the AI2D-RST corpus. The coloured bars to the right of the graphs map the colours to the values of probability density function for each diagram element type. Created using *matplotlib* 3.3.0 (Hunter, 2007) and *seaborn* 0.10.1 for Python 3.8.5. (a) Cycle. (b) Illustration. (c) Cross-section. (d) Cut-out.

We first convert each blob to greyscale. Each pixel in a greyscale image is represented by a value between 0 and 255, which encodes its brightness: 0 stands for black, whereas 255 stands for white. We then describe the brightness of the entire blob by calculating a greyscale histogram with sixty-four bins, each of which covers a range of values. To exemplify, the first bin (out of sixty-four) covers values from 0 to 3. If the value of a pixel falls within this range, the value for the first bin increases by one. Distributing all the pixels across the sixty-four bins provides a sixty-four-dimensional vector that describes the brightness of the blob.

We then extract uniform Local Binary Patterns (LBPs; Ojala *et al.*, 1996) to represent the texture of each blob using the *scikit-image* library for Python (van der Walt *et al.*, 2014). Uniform LBP examines the neighbourhood of each pixel within a prespecified window and encodes information about that pixel

neighbourhood using binary values: if the value of a neighbouring pixel is lower than the value of the current pixel, the neighbour receives a value of 0. Conversely, if the value is larger, the neighbouring pixel receives a value of 1. Uniform LBP collects this information into a vector of zeros and ones. This vector is then quantified by counting the number of transitions from 0 to 1 and 1 to 0. The number of transitions is aggregated into a histogram to describe the distribution of binary patterns in the image. We use LBP to examine the twenty-four neighbouring pixels positioned within a radius of three pixels from the centre pixel; this provides a twenty-six-dimensional vector for each blob. Finally, we concatenate the sixty-four-dimensional vector for brightness and the twenty-six-dimensional vector for texture into a ninety-dimensional feature vector. These features were extracted for each blob in the AI2D dataset.

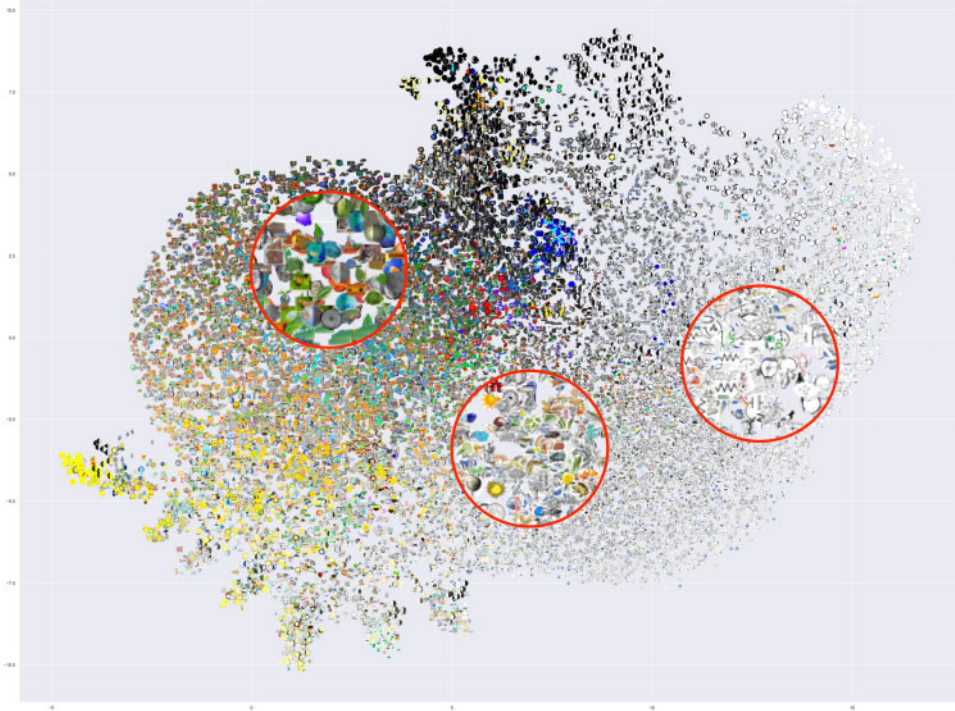


Fig. 10 Mapping the ninety-dimensional feature space for grey histogram and LBPs to two dimensions for plotting using the UMAP dimensionality reduction algorithm. Each blob is represented by its thumbnail. The three loupes demarcated in red zoom into different regions of the plot. Created using the *matplotlib* 3.3.0 (Hunter, 2007) and *seaborn* 0.11.0 libraries for Python 3.8.5.

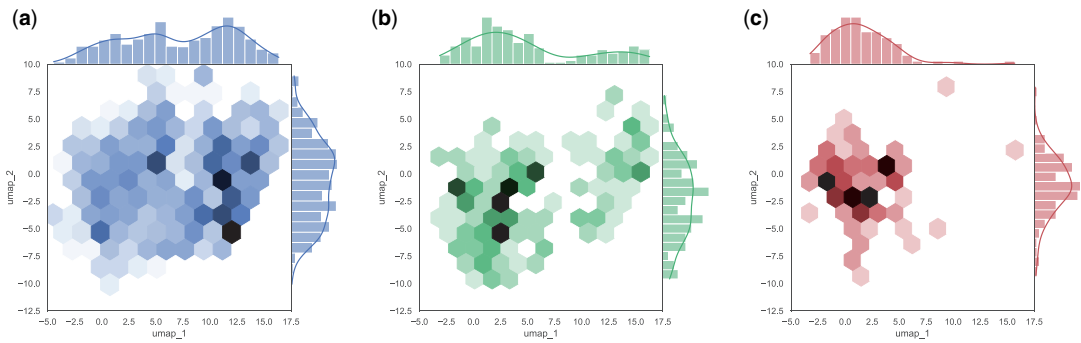


Fig. 11 The distribution of blob centroids across three structural categories in AI2D-RST. The marginal plots show histograms with 20 bins and a KDE for each UMAP dimension. Created using the *matplotlib* 3.3.0 (Hunter, 2007) and *seaborn* 0.11.0 libraries for Python 3.8.5. (a) Illustration. (b) Cross-section. (c) Cut-out.

monochrome drawings, which are both suitable for depicting the internal structure of objects from a side view. *Illustrations*, in turn, are more flexible in

terms of the choice of expressive resources that they may draw on for depiction. These results are generally aligned with the findings of Dimopoulos *et al.* (2003,

text, it would be natural to ask whether the text refers to a particular entity in the diagram. This would allow aspects of the inherently dynamic nature of discourse interpretations to be injected into the annotation tasks.

Finally, we argue that our analysis shows that pre-theoretical distinctions between ‘text’ and ‘image’ rarely hold in multimodal communication, and they are seriously under-differentiating for the digital humanities or any other field concerned with multimodality (Bateman, 2014; Bateman *et al.*, 2017). While we fully agree with Arnold and Tilton (2019) on the need to develop methods that enable the large-scale analysis of various media, we have demonstrated here that this effort must be supported by a solid theoretical foundation that reveals rather than hides the complexity of multimodal communication (Bateman, 2017; Hiippala, 2021). This foundation is needed for tackling issues that are traditionally of concern to the humanities, such as trajectories of change over time, whose computational analysis is still largely limited to linguistic material.

5 Conclusion

In this article, we argued that contemporary theories of multimodality can inform computational approaches to studying multimodal communication in the field of digital humanities and beyond. By analysing two multimodal corpora consisting of primary school science diagrams, we showed how multimodality theory can reveal descriptive shortcomings that lead to analytical blind spots, which become increasingly pronounced when using computational methods as advocated by distant viewing (Arnold and Tilton, 2019). However, Arnold and Tilton (2019, p. i13) acknowledge that the framework of distant viewing cannot specify what kinds of metadata schemes are needed to describe the data in a way that support the optimal use of computational methods, but argue that developing such metadata schemes should constitute a major area of research in digital humanities.

We propose that any effort to define metadata schemes for visual and multimodal materials can be informed by multimodality theory, which allows for a semiotically appropriate treatment of diverse media and the semiotic modes they deploy. This calls for increased attention to the communicative goals set

for the artefact or situation under analysis, as these determine the extent to which individual semiotic modes must be decomposed into analytical units to achieve a sufficiently coherent description of multimodal discourse. Producing such descriptions for a wide range of historical and contemporary media will require a large-scale effort, but at the same time, grounding the analysis in semiotics offers a far stronger basis for addressing research questions that are traditionally of concern to humanities, while continuing to leverage the power of computational methods for finding patterns in large volumes of data.

Notes

1. <https://www.mturk.com>.
2. The code used for analysis is available at <https://doi.org/10.5281/zenodo.4761066>.

References

- Alikhani, M. and Stone, M. (2018). Arrows are the verbs of diagrams. In *Proceedings of the 27th International Conference on Computational Linguistics*, Santa Fe, New Mexico, USA, pp. 3552–63.
- André, E. and Rist, T. (1995). Generating coherent presentations employing textual and visual material. *Artificial Intelligence Review*, 9: 147–65.
- Arnold, T. and Tilton, L. (2019). Distant viewing: analyzing large visual corpora. *Digital Scholarship in the Humanities*, 34(Supplement 1): i3–i16.
- Bateman, J. A. (2008). *Multimodality and Genre: A Foundation for the Systematic Analysis of Multimodal Documents*. London: Palgrave Macmillan.
- Bateman, J. A. (2011). The decomposability of semiotic modes. In O’Halloran, K. L. and Smith, B. A. (eds), *Multimodal Studies: Multiple Approaches and Domains*. London: Routledge, pp. 17–38.
- Bateman, J. A. (2014). *Text and Image: A Critical Introduction to the Visual/Verbal Divide*. London and New York: Routledge.
- Bateman, J. A. (2017). Multimodale semiotik und die theoretischen grundlagen der digital humanities. *Zeitschrift für Semiotik*, 39(1–2): 11–50.
- Bateman, J. A. (2020). The foundational role of discourse semantics beyond language. In Zappavigna, M. and

- Kress, G. and van Leeuwen, T.** (1996). *Reading Images: The Grammar of Visual Design*. London: Routledge.
- Kress, G. and van Leeuwen, T.** (2001). *Multimodal Discourse: The Modes and Media of Contemporary Communication*. London: Arnold.
- Lang, S. and Ommer, B.** (2018). Attesting similarity: Supporting the organization and study of art image collections with computer vision. *Digital Scholarship in the Humanities*, **33**(4): 845–56.
- Lenke, J. L.** (2005). Multimedia genres and traversals. *Folia Linguistica*, **39**(1–2): 45–56.
- Mann, W. C. and Thompson, S. A.** (1988). Rhetorical structure theory: Toward a functional theory of text organization. *Text*, **8**(3): 243–81.
- McInnes, L., Healy, J., Saul, N. and Grossberger, L.** (2018). UMAP: Uniform manifold approximation and projection. *The Journal of Open Source Software*, **3**(29): 861.
- Mondada, L.** (2019). Contemporary issues in conversation analysis: Embodiment and materiality, multimodality and multisensoriality in social interaction. *Journal of Pragmatics*, **145**: 47–62.
- Münster, S. and Terras, M.** (2020). The visual side of digital humanities: A survey on topics, researchers, and epistemic cultures. *Digital Scholarship in the Humanities*, **35**(2): 366–89.
- O'Halloran, K. L., Tan, S., Pham, D.-S., Bateman, J. A. and Vande Moere, A.** (2018). A digital mixed methods research design: Integrating multimodal analysis with data mining and information visualization for big data analytics. *Journal of Mixed Methods Research*, **12**(1): 11–30.
- Ojala, T., Pietikäinen, M. and Harwood, D.** (1996). A comparative study of texture measures with classification based on featured distributions. *Pattern Recognition*, **29**(1): 51–9.
- Oviatt, S. and Cohen, P. R.** (2015). *The Paradigm Shift to Multimodality in Contemporary Computer Interfaces*. San Rafael (CA): Morgan & Claypool.
- Purchase, H. C.** (2014). Twelve year of diagrams research. *Journal of Visual Languages and Computing*, **25**(2): 57–75.
- Richards, C.** (2017). Technical and scientific illustration: picturing the invisible. In Black, A., Luna, P., Lund, O. and Walker, S. (eds), *Information Design: Research and Practice*. London: Routledge. pp. 85–106.
- Smits, T. and Ros, R.** (2020). Quantifying iconicity in 940k online circulations of 26 iconic photographs. In *Proceedings of the Workshop on Computational Humanities Research (CHR 2020)*. pp. 375–384.
- Steen, F. F., Hougaard, A., Joo, J., et al.** (2018). Toward an infrastructure for data-driven multimodal communication research. *Linguistics Vanguard*, **1**: 20170041.
- Stöckl, H.** (2020). Linguistic multimodality—multimodal linguistics: a state-of-the-art sketch. In Wildfeuer, J., Pflaeging, J., Bateman, J. A., Seizov, O. and Tseng, C. (eds), *Multimodality: Disciplinary Thoughts and the Challenge of Diversity*. Berlin: De Gruyter, pp. 41–68.
- Sulubacak, U., Caglayan, O., Grönroos, S.-A., et al.** (2020). Multimodal machine translation through visuals and speech. *Machine Translation*, **34**(2–3): 97–147.
- Svensson, P.** (2010). The landscape of digital humanities. *Digital Humanities Quarterly*, **4**(1).
- Taboada, M. and Habel, C.** (2013). Rhetorical relations in multimodal documents. *Discourse Studies*, **15**(1): 65–89.
- Thomas, M.** (2014). Evidence and circularity in multimodal discourse analysis. *Visual Communication*, **13**(2): 163–89.
- Thomas, M.** (2020a). Making a virtue of material values: tactical and strategic benefits for scaling multimodal analysis. In Wildfeuer, J., Pflaeging, J., Bateman, J. A., Seizov, O. and Tseng, C. (eds), *Multimodality: Disciplinary Thoughts and the Challenge of Diversity*. Berlin: De Gruyter, pp. 69–91.
- Thomas, M.** (2020b). Multimodality and media archaeology: Complementary optics for looking at digital stuff?. *Digital Scholarship in the Humanities* (doi: 10.1093/llc/fqaa024).
- Tversky, B., Zacks, J., Lee, P. and Heiser, J.** (2000). *Diagrams 2000: theory and application of diagrams*. In: *Lines, Blobs, Crosses and Arrows: Diagrammatic Communication with Schematic Figures*. Berlin: Springer, pp. 221–230.
- van der Walt, S., Schönberger, J. L., Nunez-Iglesias, J., et al.** (2014). scikit-image: image processing in Python. *PeerJ*, **2**(e453).
- Waller, R. H. W.** (2012). Graphic literacies for a digital age: The survival of layout. *The Information Society*, **28**(4): 236–52.
- Waller, R. H. W.** (2017). Practice-based perspectives on multimodal documents: Corpora vs connoisseurship. *Discourse, Context & Media*, **20**: 175–90.
- Ware, C.** (2012). *Information Visualization: Perception for Design*. 3rd edn. Amsterdam: Elsevier.
- Watanabe, Y. and Nagao, M.** (1998). Diagram understanding using integration of layout information and

- textual information. In *Proceedings of the 36th Annual Meeting of the Association for Computational Linguistics and 17th International Conference on Computational Linguistics (ACL'98/COLING'98)*. Association for Computational Linguistics, Montreal, Quebec, Canada, pp. 1374–80.
- Wevers, M. and Smits, T.** (2020). The visual digital turn: using neural networks to study historical images. *Digital Scholarship in the Humanities*, 35(1): 194–207.
- Wildfeuer, J., Pflaeging, J., Bateman, J. A., Seizov, O. and Tseng, C. (eds)** (2020). *Multimodality: Disciplinary Thoughts and the Challenge of Diversity*. Berlin: De Gruyter.