The Big Picture: A Visual Exploration of the Reciprocal Image of Italy and China through Search Engines

Giulio Fagiolini Politecnico di Milano Via Durando 10 20158 Milano +39 0223997813 giulio.fagiolini@mail.polimi.it Paolo Ciuccarelli Politecnico di Milano Via Durando 10 20158 Milano +39 0223997813 paolo.ciuccarelli@polimi.it YANG Lei Chinese Museum of Digital Art Fuxing Road 9A Haidian District, Beijing +86 1059802310 yanglei@modachina.org

ABSTRACT

This experiment consists in the collection, categorisation and visualisation of 4,800 images from the reciprocal national internet domains of Italy and China. The digital world is here considered not only in terms of the impact of new technologies on social life, but also as a resource for the real world as a political and social space[11]. In a context where the language barrier presents a big obstacle, images can be a medium for cultural analysis by exploiting both their visual properties and their intrinsic storytelling capabilities[8]. Thanks to today's massive data production, we are now able to conduct analyses that were not possible before. This experiment is an attempt to investigate how two radically different cultures see each other through the images collected on the web. The visual characteristics of these artefacts, together with the uses to which they are put, provide a valuable tool for the investigator. We hope that this work will provide insight into the big picture for the general reader while offering the specialist a practical tool to test hypotheses and intuitions.

Categories and Subject Descriptors

H.5.m [Information Interfaces and Presentation]: Miscellaneous.

H.3.7 [Information Storage and Retrieval]: Digital Libraries – *Collection*.

J.4.7 [Social and Behavioural Sciences]: Sociology.

General Terms

Measurement, Documentation, Design, Experimentation.

Keywords

Digital Methods, Data Visualization, Cultural Analytics.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

1. INTRODUCTION

The aim of the project was to examine the peculiarities of the narrative of both countries in one another's web space. The exponential growth of non-professional and professional media producers has created a new cultural situation as well as a challenge to our normal ways of tracking and studying culture[3]. Thanks to this massive production of data we were able to make a number of analyses that were not possible previously. The questions we were interested in were, first, whether we could use the collection of images found in the reciprocal web of Italy and China as a tool to investigate the perception of respective national identities, and, second, what kind of insights these images would provide.

2. MEANING OF IDENTITY

The word "identity" envelops manifold meanings. We will use "national identity" specifically to refer to the concept of type identity: labels applied to persons (and in this case nations) who share or are thought to share some characteristics in appearance, behavioural traits, beliefs, attitudes, values, knowledge, opinions, experience, historical commonalities, and so on. Since the concept of "national identity" also acts as a social category[1], it is subject to the products of human thinking, discourse, and action. This makes the meaning of "national identity" variable over time and location, and therefore it follows that the result of our analysis is valid strictly for the moment in which it was made and for the place in which it was made. Dialogic processes can serve to determine such identities (as do forms of communication and interaction between nations) by giving others a method of defining a perceived national identity. Understanding communication as a dialogic process allows us to extract information on the subjects involved, the perceiver and the perceived.

3. CHINA AND ITALY

China has always been a vast territory with a strong cultural identity, historically powered by the heritage of written tradition[6]. Chinese culture has developed independently from the west for thousands of years. Despite China's long history of strong trade relations, it is only after the reforms carried out under Deng Xiaoping in the 1980s that China began opening her doors to the rest of the world[6]. In order to manage these transitions and changes without losing its role, the government exerts strict control over the media and the web. The result is an incredibly complex society that, while extremely dynamic and futuristic, remains essentially separated from the rest of the world. Bearing this in mind, what really makes China an interesting country for this experiment is the cultural barrier created by the language.

On the other side we have Italy, a country with a media landscape similar to many of the non-English-speaking western countries.

The advent of digital technologies in journalism is still recent, and the Italian media is in the middle of a process of evolution and redefinition. The general movement toward the web has resulted in a loss of power and revenue for the established media sources (so-called generalists), and has obliged them to search for a new form, while new smaller publishers are also trying to find their own space.

4. METHODOLOGY

The background to this research combines two approaches developed by the Digital Methods Initiative of Amsterdam and the Software Studies Initiative of New York. The first method, which considers the digital sphere both as a measure of the impact of new technologies on the user and as a resource used by the real world as a political and social space[11], introduces the term "online groundedness" in an effort to conceptualise the research that follows the medium, to capture its dynamics and make grounded claims about cultural and societal change[9]. The second approach focuses on research into software and the way computational methods can be used for the analysis of massive data sets and data flows in order to analyse large collections of images. "If media are 'tools for thought' through which we think and communicate the results of our thinking to others, it is logical that we would want to use the tools to let us think verbally, visually, and spatially."[4] Other works share a similar approach with this project. In particular the narrative structure that characterises the work is matched by the Selfiecity project (http://selfiecity.net/), where a similar analysis is employed to investigate selfies taken around the world.

4.1 Selection of sources

Having decided to examine the perceived identities of these nations in their mutual web-spaces through images and to pay close attention to how this identity is "broadcasted", search engines, being a crucial point of entrance and exploration of the web, seemed a natural place to start. The two main sources for the collection of data were therefore the two main image-search engines of the two countries. Google's position as the main search engine in Italy (we refer here specifically to the national domain google.it), is mirrored by Baidu in China, which commands about two-thirds of the booming search market there[7]. To add a further layer to the research, we employed Google's advanced search instruments to conduct a second series of queries limited to a selection of domains concerning specific news websites that carried particular meaning for either country. Thus the collection included 2,400 images for each data set obtained by searching for the translated name of one nation in the local nation's web space: 900 images retrieved directly from the respective search engine and 300 from five different news websites scraped via the search engine.

4.2 Data collection

In order to ensure that research on the images was as objective as possible, it was crucial to isolate it from personal computer and search engine use. Some rules were implemented for this purpose: • Log out from any Google service

• Delete all customisation and localization services related to social networks and browser history

• Empty the search engine's cache

Because data collection from the Chinese web was done in mainland China, it was not necessary to use proxy or other software to simulate the originating location of the queries. Each query was conducted from the country of the specific domain. The collection of images was carried out between 01-15/02/2013 for images pertaining to China, and between 01-15/03/2013 for images regarding Italy. The period in question is fundamental for the analysis of the content. The results show a combination of collective memories, everyday narratives and the peculiarities of each day: a sampling of separate moments, seasons, amplifications and contractions of time as they appeared at the instant in which they were harvested.

4.3 Data processing

Before beginning to visualise, it was necessary to understand all the data enclosed in the images. We first measured the properties in each image by using the QTIP digital image processing application that provided us with measurement files listing the mean values of brightness and the use of red, green and blue in each image. Then, to provide a qualitative dimension to the research, the images selected were manually categorised. They were organised into a hierarchical and multiple taxonomy. This allowed us to track the characteristics of each image and identify the main thematic clusters. We ended up with around 100 sub-categories belonging to seven main categories: Architecture, Disaster report, Economics, Nature, Non-photo, Politics, Society, and Sport.

5. VISUALISATION

The visualisation element is divided into two parts: we will first discuss the production of some of the actual visualisations, then the development of the container.

5.1 Visualising visual features

The first intention was to take a step back and compare the images of the two datasets in terms of their visual features. We relied on the Cultural Analytics tools and techniques developed by the Software Studies Initiative at the University of California, San Diego[2]. By exploring large image sets in relation to multiple visual dimensions and using high resolution visualisations, the Cultural Analytics approach allows us to detect patterns which are not visible with standard interfaces for media viewing. In contrast with standard media visualisations which represent data as points, lines, and other graphical primitives, Cultural Analytics visualisations show all the images in a composition[5].

We therefore used ImageMontage and ImageSlice to produce several large canvases of thumbnails showing the respective visual features characterising the different collections.



Figure 1. Example of comparison between the ImageMontages of the images about China (left) and Italy (right) ordered by saturation.



Figure 3. Example of comparison between the ImageSlices of the images about China (top) and Italy (bottom) ordered by hue.

These representations allow us to identify easily the points of continuity and discontinuity between the visual features of the two data sets, while selective ImageMontages (Figure 3) quantify the differences according to each step of the value.



Figure 2. Comparison between the selective ImageMontage of the images about China (left) and Italy (right) ordered by steps of hue.

As we can see from the visualisations, each nation has a specific Local Colour: visual attributes and dominant tones, which relate to specific cultural territories[2].

5.2 Visualising categories

A specific visual model was developed to visualise the categories and its subcategories. It shows the main category as the central bubble around which the sub-keywords are disposed in circles for the identification of relevant issues. Each image is tagged with one or more keywords/sub-keywords, and the dimension of each bubble is proportional to the number of images tagged with a keyword or sub-keyword.



Figure 4. Bubble packing of the subcategories related to the category Architecture from the database about China



Figure 5. Bubble packing of the subcategories related to the category Architecture from the database about Italy

Although the visual model does not show which keywords have images in common, it highlights the differences between the two databases and therefore between the two countries.

In order to compare the relevance of each keyword to each of the sources, we made a series of bar charts. Each one represents the profile of a single source. In this way we could easily contrast the different "vocations" of the sources by highlighting the space given to each topic.



Figure 6. Bar charts comparing the profile of two of the sources from Italy.

Figure 6 shows the example of a comparison between the profiles of the newspapers *Corriere della Sera* and *Fatto Quotidiano*. In this case it is clear how the latter stresses economics and politics while *Corriere* has a relatively balanced profile.

5.3 Showcase

The culmination of our experimental project has been the creation and development of the website http://thebigpictu.re where the main visualisations have been collected. In the process of creating this interface our focus has remained on the same idea from which this project originated: to increase awareness of the way we see and the way we are seen by a culture radically different from our own. This was done by making a tool which makes the topic comprehensible to outsiders, without the need for simplification, as well as to specialists in the field.

From a data visualisation point of view, the biggest challenge was to find an appropriate structure: simplified enough to show the big picture emerging from the data and detailed enough to preserve all the interesting details in the data. We acted on this in two ways: first, we decided to set up the narration consistently on a comparative level; and second, to give the user a tool for a multifaceted exploration of data. Keeping the visualisation and the storytelling on a comparative level helped to keep the exploration clean and structured, which also enabled us to explain each level of the research. The exploration tool, as a personal tool for navigating the data set, represents the last phase of the work.

The exploration tool makes use of 'Elastic lists', a user interface component for facet browsers developed by Moritz Stefaner. It aims to enrich current interfaces with additional visual cues about the relative weights of metadata values, as well as how that weight differs from the global metadata distribution[9].

6. CONCLUSIONS

The method we developed allows us to draw some conclusions about the dominant visual features and the dominant topics characterising the two nations in the other nation's webspace.

Regarding the visual features, the main discontinuity between the two nations pertains to the median hue and the median saturation of the pictures. The data set of images about China is strongly characterised by the colour red which has little presence in the data set about Italy, which conversely, is represented by the colour blue. This particularity is experienced also through the analysis of the saturation. In terms of the topics, there are many other strong points of discontinuity. In particular, it is interesting to note how the images of the architectural heritage of the two nations differs (Figure 4 and 5). The representation of Italy is pretty diverse, with traditional architecture from almost everywhere in the country and natural landscapes from the coast, mountains and lakes. What is missing is a representation of modernity. China, on the other hand, is represented by a very narrow range of attractions that gather all the attention. Much space is given to the Great Wall, and this is almost the only recognisable traditional structure within the data set (it shares the same presence with Venice, without even being a city). We also found a strong presence of modern architecture, with many views of skyscrapers and infrastructures, even even if some had a cynical and ironic disposition.

To conclude, we can say that the work allows the user not only to explore all the singular elements of the database but also to focus on the database as a whole. As the title states, the overall purpose and outcome is to show a big picture including all the facets that make it unique.

6.1 Limits

The method used to conduct this analysis had some inherent limitations.

There are many variables that influence the results: the observer,

the place from which s/he is observing, the moment in which s/he is observing, the duration of the observation, the "focusing", and the distance between the observer and the observed. All these considerations can also be referred to reality, i.e. they are also true of the non-virtual world: if we asked a heterogeneous group of people to describe a certain place, we would end up with as many subjective descriptions of the reality as there were people in the group. However, the categorisation of 4,800 images was carried out manually by one person, in order to limit changes to the parameters of evaluation in the process.

6.2 Developments

Future research might benefit from the use of automated content analysis technology. By automating the categorisation in the analysis, the tagging process could be made more neutral. The same methodology used and developed during this project could also be applied to different data sets. For example, it would be interesting to conduct the same experiment on one of the newspaper archives, where the existing tagging system would make the process of categorisation easier and more automated.

7. REFERENCES

- [1] Fearon, J. *What is Identity (as we now use the word)*?, Unpublished paper. Available online. Stanford University, 1999.
- [2] Hochman, N. and Schwartz, R. *Visualizing Instagram: Tracing Cultural Visual Rhythms*. Technical Report. Association for the Advancement of Artificial Intelligence, 2012.
- [3] Manovich, L. *Cultural Analytics: Visualing cultural patterns in the era of more media*. Domus, New York, 2009.
- [4] Manovich, L. Software Takes Command. Bloomsbury Academic, London, New York, 2013.
- [5] Manovich, L. Media Visualization: Visual Techniques for Exploring Large Media Collections. Media Studies Futures, ed. Kelly Gates. Blackwell, 2012.
- [6] Pietrasanta, B. *L'ideogramma al Neon*. Lupetti, Milano, 2009.
- [7] Reinsberg, R. *Baidu*. MIT Technology Review http://technologyreview.com/article/416835/baidu/, 2009.
- [8] Ricci, D. Seeing what they are Saying: diagrams for sociotechnical controversies. Doctoral thesis. Politecnico di Milano, 2010.
- [9] Rogers, R. *The End of the Virtual: Digital Methods*. Text prepared for the Inaugural Speech, Chair, New Media & Digital Culture, University of Amsterdam, 2009.
- [10] Stefaner, M. and Muller, B.*Elastic lists for facet browsers*. University of Applied Sciences, Potsdam, 2007.
- [11] Weltevrede, E. *Thinking Nationally with the Web.* M.Sc. dissertation., University of Amsterdam, 2009.