# Design Report:
# Digital Methods for Large-Scale Film Analysis

## 1. Films as Visual Corpora

In digital humanities (DH), scholars have extensively explored automatic research methods for identifying latent patterns in large corpora of literary texts. For instance, the Culturomics project of Jean-Baptiste Michel *et al.* studied historical trends in language usage based on a corpus of more than 5 million books, which was believed to cover approximately 4% of all books published worldwide.[1] With a focus on periodicals from the United Kingdom (UK), Thomas Lansdall-Welfare *et al.* applied named entity detection to understand gender representation in 35.9 million news articles that represented 14% of all newspapers published in the UK between 1800 and 1950.[2] In these studies, corpus size and coverage are often attached with great importance.

In the vast world of visual culture, however, the establishment of a digital methodology for interpreting large-scale visual corpora is still in its infancy. The coinage of *Distant Viewing* by Taylor Arnold and Lauren Tilton in 2019 has provided a theoretical framework for digitally analysing visual texts and raised scholarly attention to the potential use of Computer Vision (CV) techniques in DH. Nonetheless, research on massive visual corpora is still limited by two drawbacks.[3] First, scholars have shown much interest in databases of still images, but moving images, such as films, music videos, and television shows, still receive relatively less attention.[4] Second, studies investigating moving images on large scales often focus on a limited range of selected texts, such as television series or film franchises, to reduce the workload of collecting

---

[1] Jean-Baptiste Michel and others, 'Quantitative Analysis of Culture Using Millions of Digitized Books', *Science*, 331.6014 (2011), 176–82.

[2] Thomas Lansdall-Welfare and others, 'Content Analysis of 150 Years of British Periodicals', *Proceedings of the National Academy of Sciences*, 114.4 (2017), E457–65.

[3] Taylor Arnold and Lauren Tilton, 'Distant Viewing: Analyzing Large Visual Corpora', *Digital Scholarship in the Humanities*, 34.Supplement_1 (2019), i3–16.

[4] E.g. Ana Jofre and others, 'Faces Extracted from Time Magazine 1923-2014', *Journal of Cultural Analytics*, 5.1 (2020); Taylor Arnold, Lauren Tilton, and Justin Wigard, 'Automatic Identification and Classification of Portraits in a Corpus of Historical Photographs' (presented at the Computational Humanities Research Conference 2022, Antwerp, Belgium: University of Antwerp, 2022), pp. 25–35.

metadata.[5] Although this approach allows researchers to build corpora of moving images more efficiently, it also sets a methodological limitation for their research scope.

With the knowledge of these two methodological drawbacks, my project explored the application of large-scale film corpora in DH. This project drew on the method of building large-scale corpora of literary texts, where the research scope was set as films released in a specific temporal period and region instead of a franchise. The subject of inquiry in my project is the film industry of mainland China between 1949 and 1966 (also known as the "Seventeen Years Period" in Chinese cinema studies). This subject is considered a topic that typically lacks quantitative studies. Researchers in both Anglophone and Chinese-language scholarship have reported that the study of cinema history in Maoist China (1949-1976) excessively relied on the close reading of selected film texts, which might limit the representability of the findings from existing research.[6] This research gap provided a practical context for applying digital methods to complement previous studies. Based on the research practice of this project, a dialogue between film archiving and computer vision unfolds.

## 2. Building the Corpus

### 2.1. Film Archiving in China

Collecting metadata for a large-scale film corpus requires rigorous use of documentation published by film archives. In a study on the Soviet film director Dziga Vertov, Adelheid Heftberger demonstrated the benefits of using documents provided by official film archiving institutions to understand the conditions of film production.[7] However, it was revealed in this

---

[5] Taylor Arnold, Lauren Tilton, and Annie Berke, 'Visual Style in Two Network Era Sitcoms', *Journal of Cultural Analytics*, 4.2 (2019); Alina El-Keilany, Thomas Schmidt, and Christian Wolff, 'Distant Viewing of the Harry Potter Movies via Computer Vision' (presented at the The 6th Digital Humanities in the Nordic and Baltic Countries Conference, Uppsala, Sweden: Uppsala University, 2022), pp. 33–49.

[6] Binyan Pi, '中国"十七年"电影研究的回顾与反思 [A Review and Reflection of Research on Chinese Cinema in the "Seventeen Years Period"]', *Bohai University Journal (Social Science Volume)*, 39.2 (2017), 100–104; Zhuoyi Wang, *Revolutionary Cycles in Chinese Cinema, 1951-1979*, 1st edn (New York: Palgrave Macmillan, 2014), pp. 11–14.

[7] Adelheid Heftberger, *Digital Humanities and Film Studies: Visualising Dziga Vertov's Work*, Quantitative Methods in the Humanities and Social Sciences (Cham: Springer International Publishing, 2018).

project that reliance on official documents as the only data source would potentially impair the integrity of the corpus.

Initially, I focused on *The Catalogue of Chinese Artistic Films* (shortened as *The Catalogue*) and *The Encyclopedia of Chinese Films* (*The Encyclopedia*) as the main sources for data collection in this project.[8] These two books were published by the China Film Archive (CFA), the largest film preservation agency in the People's Republic of China (PRC). Both books were supposed to be authoritative publications, but there was a significant discrepancy in their coverage of films produced by private film companies in Shanghai between 1949 and 1953.

| Source | All Entries | Films from Private Studios in Shanghai |
|---|---|---|
| *The Catalogue* | 656 | 24 |
| *The Encyclopedia* | 722 | 85 |
| This Project | 732 | 86 |

*Table 1*. Entry counts in the two books and this project.[9]

The difference between the two books raised caution regarding the use of official documentation. *The Catalogue* published in 1981 omitted 61 films produced by private companies in Shanghai that could be found in *The Encyclopedia* published in 2001 (*Table 1*).[10] In addition to the efforts paid by CFA researchers in discovering lost media, this difference is also related to the relationship between mainland China and Taiwan. Many private filmmaking institutions in Shanghai had intimate relationships with the government of the Republic of

---

[8] Fenglan Jin and Zhengduo Yang, 中国艺术影片编目 *[The Catalogue of Chinese Artistic Films]* (Beijing: Culture and Art Publishing House, 1981); China Film Archive, 中国影片大典：故事片，舞台艺术片，

*1949.10-1976 [Encyclopaedia of Chinese Films, 1949.10-1976]*, 1st edn (Beijing: China Film Press, 2001).

[9] *The Encyclopedia* included 27 entries produced in Hong Kong or Taiwan. These films are excluded in the statistics.

[10] According to a 2005 report of China Film Archive, there were 2298 films produced in mainland China between 1905 and 1949, but only 331 of them were preserved. See Tongsheng Zhao, 电影档案概论

*[Introduction to Film Archiving]* (Beijing: China Radio, Film and Television Press, 2005).

China (ROC) in the 1940s.[11] After the PRC took over mainland China in 1949 and the ROC government moved to Taiwan, the tension between the mainland and Taiwan lasted for decades until it started to ease gradually in 1979.[12] This tension was also reflected in *The Catalogue*'s selective inclusion of films from Shanghai's private film companies, as these companies used to be cultural institutions serving the ROC. The representability of filmographies published by official organisations could be directly influenced by political demands, so the combination of official and informal sources would be necessary to ensure the coverage of the film corpus. By combining and comparing the information from *The Catalogue*, *The Encyclopedia*, the filmography of mainland China written by Donald J. Marion,[13] and two online databases (*1905.com* and *Douban*), I constructed a collection containing the metadata for 732 films as the base corpus for this project.[14]

## 2.2. Data Modelling

I selected *The Catalogue* as the main reference for data modelling. While *The Catalogue* was limited in coverage due to its selective inclusion of films from Shanghai's film companies, this 1981 book still contained valuable first-hand materials about the conditions of filmmaking in mainland China. These materials were collected by CFA-funded researchers in collaboration with state-owned film studios of the PRC and provided some exclusive information, such as the number of film reels and special techniques (e.g. widescreen film reels imported from the Soviet Union) used by a film. Therefore, I used texts extracted from *The Catalogue* with Optical Character Recognition (OCR) as the primary source for data modelling.

---

[11] Yingjin Zhang, *Chinese National Cinema*, 1st edn (New York; London: Routledge, 2004), pp. 97–99.
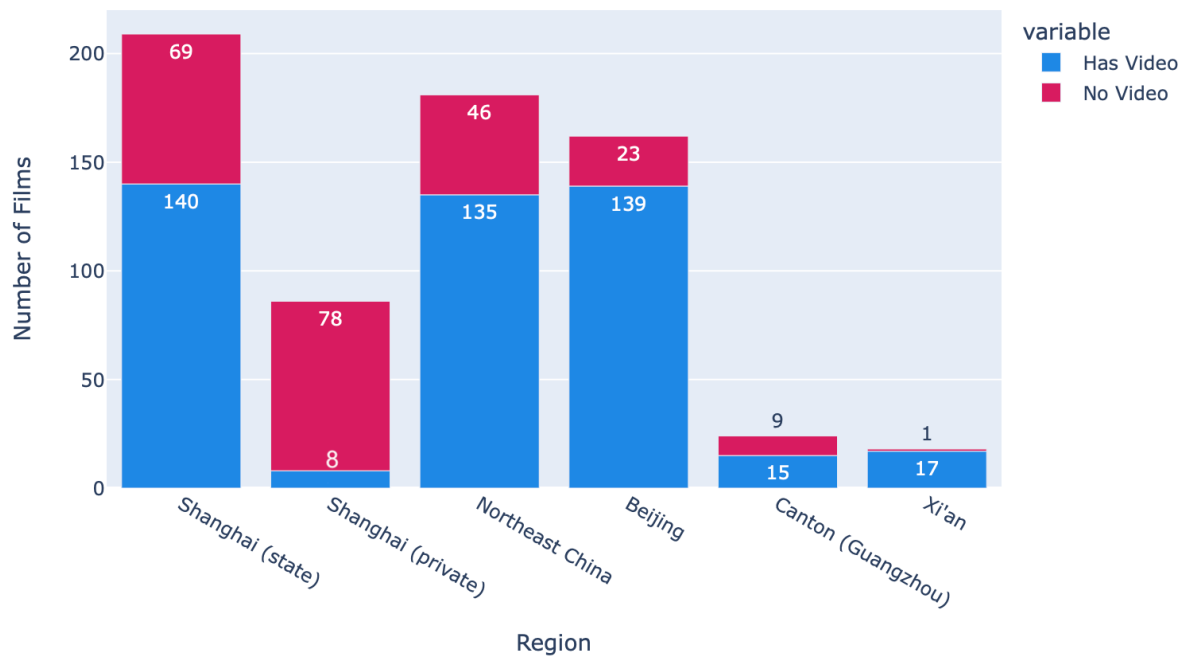
[12] Kevin G. Cai, *Cross-Taiwan Straits Relations since 1979 Policy Adjustment and Institutional Change across the Straits* (Singapore; Hackensack, NJ: World Scientific, 2011).

[13] Donald J. Marion, *The Chinese Filmography: The 2444 Feature Films Produced by Studios in the People's Republic of China from 1949 through 1995* (Jefferson, NC; London: McFarland & Co, 1997).

[14] 'cinematic', 2023 <https://github.com/El-Mundo/cinematic> [accessed 5 January 2025].

| Attributes | Available Samples |
|---|---|
| Title (Chinese) | 732 |
| Translated Title (English) | 732 |
| Release Year | 732 |
| Studio(s) | 732 |
| Colour | 731 |
| Film Reels | 656 |
| Staff and Crew | 732 |

*Table 2*. The data model.



*Figure 1*. Video data coverage.[15]

---

[15] Considering some geopolitical conventions of China, the Northeast China category is presented to include one film collaboratively made by Changchun Film Studio and Harbin Film Studio, and the Canton category includes six entries produced in collaboration with Hong Kong institutions.

I designed the data model based on the attributes available in the book (*Table 2*). Five filmmaking regions were then identified by mapping the films to the geographical locations of film institutions. After removing entries from minor filmmaking regions (where the total number of films produced between 1949 and 1966 was less than ten) and co-produced entries that could not be located in a single region, a list with 680 entries was generated, and the video data of 454 entries were collected from *1905.com*, the official film database of mainland China run by China Central Television (*Figure 1*). This video dataset covered 75.59% of all films produced by state-run studios in the five investigated regions between 1949 and 1966.

The data modelling process also revealed a critical limitation of Chinese OCR technology. Although *The Catalogue* was published in simplified Chinese, it contained some characters that are now only used in traditional Chinese. The mixed usage of simplified and traditional Chinese characters was common in publications in mainland China from the 1950s to the 1980s, as the replacement of traditional Chinese with simplified Chinese was a gradual process.[16] However, mainstream OCR libraries, such as Tesseract and EasyOCR, usually consider simplified and traditional Chinese as two distinct languages, resulting in inaccuracies in digitalising these books. This limitation of Chinese OCR reflects the general lack of digital preservation for historical documents published in mainland China. To ensure the precision of the OCR results, the extracted data was checked by human proofreaders before modelling.

Another issue disclosed in data modelling was the over-reliance on direct translation and Pinyin in existing filmographies of China. For instance, Marion's filmography used Pinyin unanimously to represent names of ethnic minorities, while Pinyin is typically used for Han-Chinese names. *The Encyclopedia* directly translated a 1958 film co-produced by Chinese and French institutions as *The Kite* rather than its original title *Cerf-volant du Bout du Monde*. The homogeneous usage of Pinyin assumed that all names written in the Chinese language were of Han-Chinese nationality but neglected China's cultural diversity and practices of transnational film production. With the awareness of this issue, I appended a translated list of all non-Chinese and ethnic-minority names to this project.[17] This list obeyed the conventions of Romanising different languages rather than solely relying on Pinyin. As some Soviet and French filmmakers

---

[16] Daniel Kane, *Chinese Language: Its History and Current Usage*, 1st edn (North Clarendon, VT: Tuttle Publishing, 2006).

[17] 'Romanised Names', 2023 <https://github.com/El-Mundo/cinematic/blob/master/OCR/JAVA/romanizing/roman-non-han_names.csv> [accessed 5 January 2025].

had collaborated with Chinese film studios in the investigated period and were included in the metadata, the translation of their names provided in this list may also interact with future research projects about the transnational filmmaking practices in the PRC.

## 2.3. Named Entity Detection

I used named entity detection to find all individuals appearing in the crew and staff information of the 732 films, but this method could not identify duplicate names. By manually checking the appearances of some common Chinese names that may cause name duplication, I found that name duplication was rare in mainland China's film industry in the studied period. The only known case of a duplicate name is "Li Lingyun (李凌云)," a name shared by a female actor and a male composer.

## 3. Computer Vision

### 3.1. Overview

While film archives provided access to film data, computer vision allowed me to automate the analysis of the collected data on a large scale. Recent progress in Computer Vision (CV) techniques has enabled substantial possibilities for new digital humanities (DH) methods to be applied in analysing large visual corpora.[18] Some research-tested CV approaches, including object classification, face recognition, and gender estimation, have already yielded novel insights for studies on visual culture.[19] However, automating the analysis of large-scale visual corpora with CV also requires careful procedures because biases might exist in the pre-trained datasets of CV models.[20] In this project, I tested four CV techniques to extract visual features from film texts on large scales: crowd size estimation, facial expression measuring, face

---

[18] Arnold and Tilton, pp. 4–5.

[19] Anssi Männistö and others, 'Automatic Image Content Extraction: Operationalizing Machine Learning in Humanistic Photographic Studies of Large Visual Archives', in *arXiv: Computer Vision and Pattern Recognition*, 2022; Jofre and others; Antoine Mazières, Telmo Menezes, and Camille Roth, 'Computational Appraisal of Gender Representativeness in Popular Movies', *Humanities and Social Sciences Communications*, 8 (2021), 137.

[20] Simone Fabbrizzi and others, 'A Survey on Bias in Visual Datasets', *Computer Vision and Image Understanding*, 223 (2022), 103552.

clustering, and gender estimation.[21] Although the latter two techniques were not used in the research data due to potential biases and computational infeasibility, they helped me to develop more critical thinking about CV techniques' application in DH.

## 3.2. Identifying Crowds



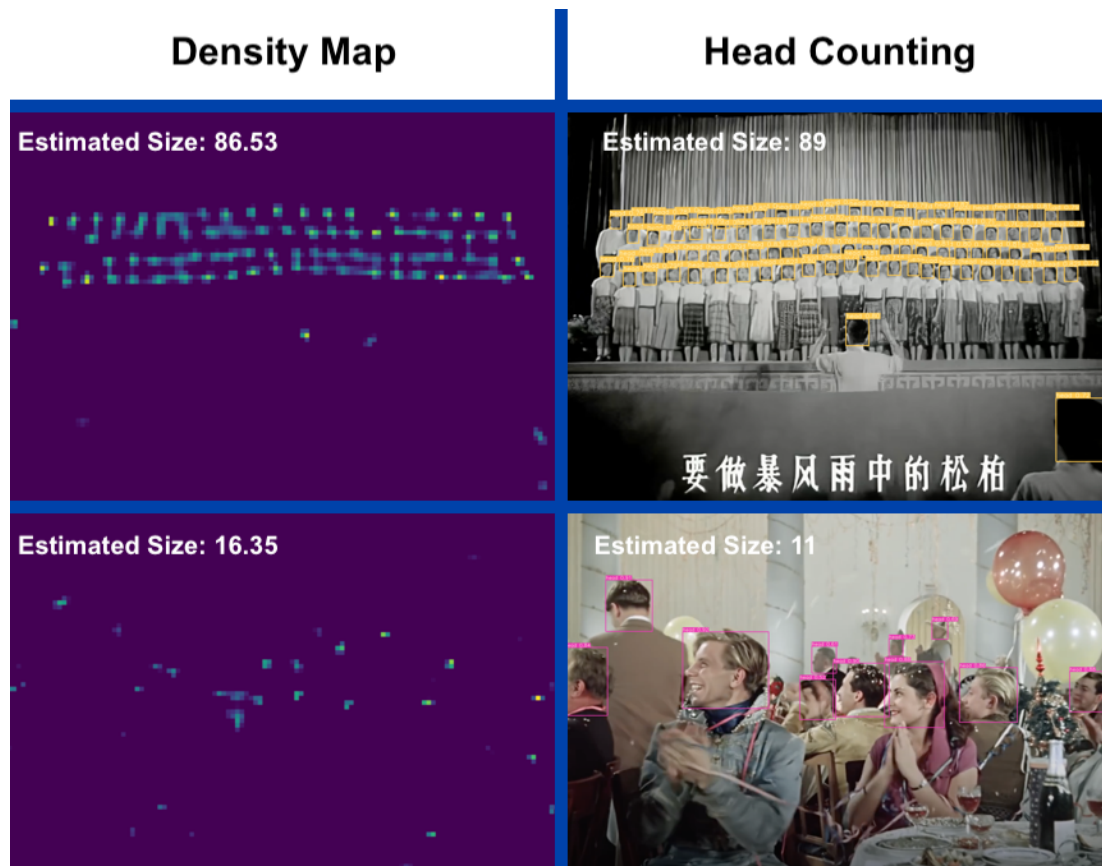*Figure 2*. A scene where a crowd is presented but not visually significant.

I combined human detection and head counting functions implemented in YOLOv8 to detect the presence and size of crowds.[22] I only considered images with more than ten detectable humans as possible crowd scenes and then applied head counting to estimate the size of the presented crowds in these images. This approach could effectively exclude scenes where a crowd was presented but lacked visual significance (*Figure 2*).

---

[21] Please see section 5 of the dissertation for the analysis of the results of crowd size estimation and facial expression measuring.

[22] Glenn Jocher, Ayush Chaurasia, and Jing Qiu, 'YOLO by Ultralytics', 2023 <https://github.com/ultralytics/ultralytics> [accessed 14 June 2023].
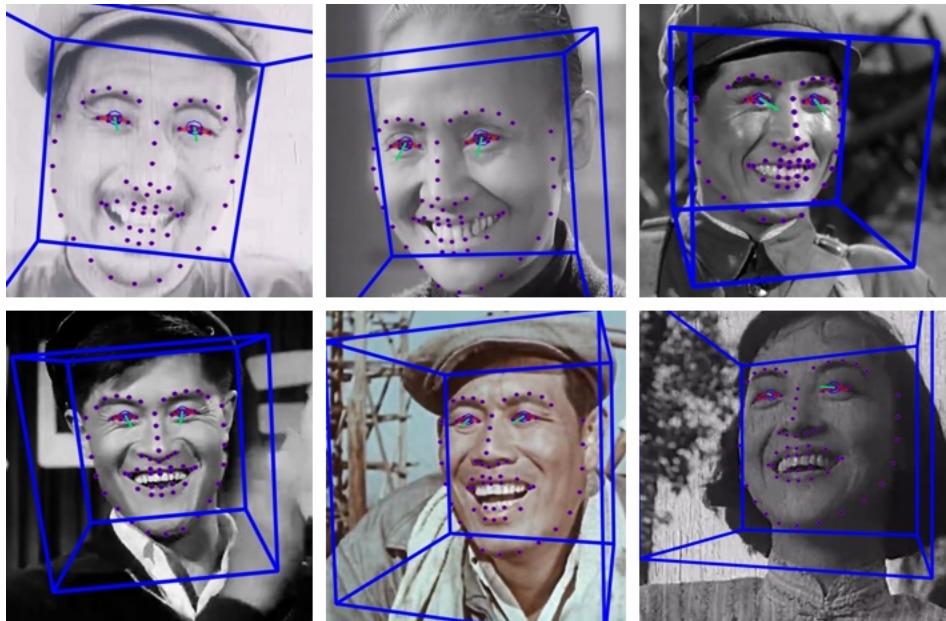
*Figure 3*. Results from the density-map algorithm and head counting.

Considering the possible low accuracy of the head-counting algorithm in dense crowds, I compared the results of head counting with another crowd-estimation algorithm based on density maps. It was found that head counting could generate more satisfactory estimations in most situations, while the density-map one could only support scenes where the characters were regularly arranged (*Figure 3*). Therefore, only the results of the head-counting algorithm were used for analysis in the research.

### 3.3. Facial Expressions



*Figure 4*. Visualisation of facial motions generated by OpenFace.

In quantitative studies of literary texts, measuring the emotional intensity of words is an effective method for identifying propaganda texts.[23] While it might be difficult to directly quantify the emotional intensity that an image conveys to its viewers, the recognition of Facial Action Units (FAU) enabled me to measure the strength of facial expressions.[24] My project implemented the FAU measuring function provided by the OpenFace library to quantify the intensity of facial expressions.[25] Because the accurate extraction of facial features required high image resolution, only faces that occupy more than one-fifth of the height of the frame were used for the detection of FAU in this project. The removal of low-resolution faces ensured the accuracy of recognition and allowed more efficient detections of faces that were significant in the visual composition (*Figure 4*).

---

[23] Travis Morris, 'Extracting and Networking Emotions in Extremist Propaganda' (presented at the 2012 European Intelligence and Security Informatics Conference, Odense, Denmark: IEEE, 2012), pp. 53–59.

[24] Seho Park and others, 'Differences in Facial Expressions between Spontaneous and Posed Smiles: Automated Method by Action Units and Three-Dimensional Facial Landmarks', *Sensors*, 20.4 (2020), 1199.

[25] Tadas Baltrusaitis and others, 'OpenFace 2.0: Facial Behavior Analysis Toolkit' (presented at the 13th IEEE International Conference on Automatic Face & Gesture Recognition, Xi'an: IEEE, 2018), pp. 59–66.

## 3.4. Face Clustering

In addition to the two methods included in my research, I performed a methodological experiment using face clustering to study the portrayal of famous film actors in different filmmaking regions. Face clustering is a CV technique that retrieves different facial images belonging to the same identity by calculating facial features as high-dimensional vectors. I experimented with this technique as a possible approach to discovering the most frequently appearing actors in films from a certain region.

Compared with face clustering, face recognition might be a more common CV method to identify characters and actors in previous studies.[26] Face recognition was highlighted by Taylor Arnold *et al.* in their comparative analysis of the narrative structures of two television series, which demonstrated the use of this technique in identifying the characters that are most actively involved in the narrative.[27] Face recognition is especially helpful for studying the presence of a predetermined set of characters. However, as films usually do not have a constant set of main characters as television series do, face clustering would be more practical than face recognition for analysing the large corpus of films. Face clustering can be performed without human supervision, which is especially meaningful in situations where the identity of the investigated actor or character cannot be determined in advance.

---

[26] Florian Schroff, Dmitry Kalenichenko, and James Philbin, 'FaceNet: A Unified Embedding for Face Recognition and Clustering' (presented at the Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, Boston, MA: IEEE, 2015), pp. 815–23.
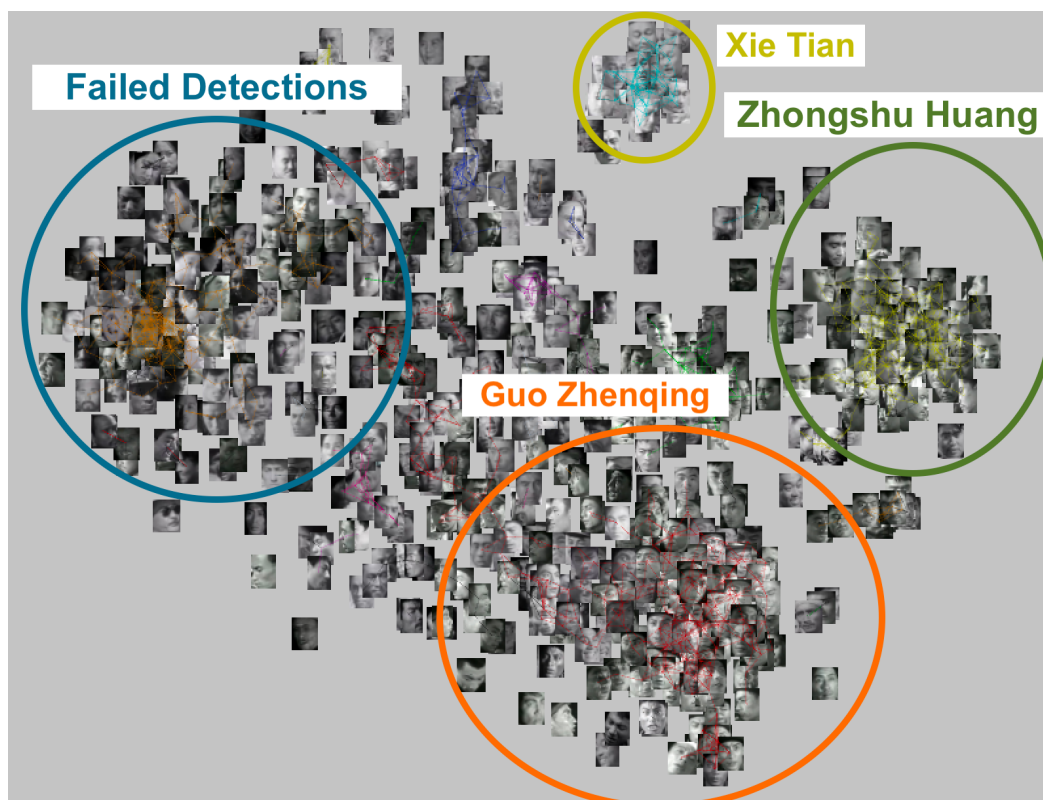
[27] Arnold, Tilton, and Berke, pp. 9–10.

*Figure 5*. Face clustering demonstration.[28]



*Figure 6*. Guo Zhenqing's face cluster.

---

[28] The face clustering application is a part of this project. See 'SimpleFaceClustering', 2023

<https://github.com/El-Mundo/SimpleFaceClustering> [accessed 5 January 2025].

To assess the feasibility of using face clustering to detect the appearances of the same actor in different films, I created an experimental dataset of 615 frames selected from three films starred by Guo Zhenqing, a famous Chinese actor in the 1950s. I used the Facenet512 model to extract the facial features and then employed the t-SNE algorithm to map the high-dimensional data of facial features onto a two-dimensional diagram.[29] As the diagram shows, this method successfully grouped Guo Zhenqing's appearances in the three films into a large cluster of face images at the diagram's bottom (*Figures 5-6*). It also recognised two actors who collaborated closely with Guo in these films. Nevertheless, it can still be noticed that a cluster of faces with random identities is shown on the left side of the diagram, which mainly contained faces found in dark environments or profile faces that could not be identified. Face clustering is a promising method for detecting unique entities in large visual corpus without human supervision, but it also has higher requirements on the resolution of images. As the corpus used in my project was historical film texts with relatively low definition, I did not apply this method on a large scale.
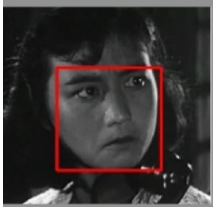
## 3.5. Gender Estimation

Gender estimation is another CV technique I practised to experiment with the large film corpus. Gender estimation has been tested by previous studies to be a useful approach for quantifying gender biases in films. For example, Ji Yoon Jang *et al.* demonstrated the stereotypical portrayals of women in 40 Korean and American films by combining this method with facial feature analysis.[30] Focusing on a larger research scope, Antoine Mazières *et al.* used gender estimation to study the evolution of women's presence in mainstream films from 1985 to 2019.[31] Although it was considered the most well-established research method among the four CV methods used in my project, the results of gender estimation were deprecated in the project due to the potential biases they may cause.

---

[29] Andrian Firmansyah, Tien Fabrianti Kusumasari, and Ekky Novriza Alam, 'Comparison of Face Recognition Accuracy of ArcFace, Facenet and Facenet512 Models on Deepface Framework', 2023, pp. 535–39.

[30] Ji Yoon Jang, Sangyoon Lee, and Byungjoo Lee, 'Quantification of Gender Representation Bias in Commercial Films Based on Image Analysis', *Proceedings of the ACM on Human-Computer Interaction*, 3.CSCW (2019), 198.

[31] Mazières, Menezes, and Roth.

| | Detected Gender | Expected Gender | Confidence |
|---|---|---|---|
| | Man | Woman | 98.6 |
| | Woman | Woman | 70.0 |
| | Man | Woman | 83.2 |
| | Man | Woman | 62.9 |

*Figure 7*. Biased gender estimations.

I applied the gender estimation model of the OpenFace library in an attempt to study gender representation in Chinese cinema.[32] However, when manually checking the results of this algorithm, I found that the result generated by this model frequently labelled women as men and caused the statistics about gender presence to be significantly biased toward men (*Figure 7*). In addition, I observed that this algorithm generally performed better in detecting gender in Chinese opera films, which seemed to result from the brighter colours and more explicit depictions of gender differences, especially the highlight of makeup, in opera films.

---

[32] Sefik Ilkin Serengil and Alper Ozpinar, 'An Evaluation of SQL and NoSQL Databases for Facial Recognition Pipelines' (Cambridge Open Engage, 2023) <https://doi.org/10.33774/coe-2023-18rcn>.

| Category | Accuracy (High Resolution) | Accuracy (Low Resolution) |
|---|---|---|
| African | | |
| With Makeup | 80% | 73.3% |
| Without Makeup | 76.6% | 70% |
| Caucasian | | |
| With Makeup | 83.3% | 70% |
| Without Makeup | 86.6% | 76.6% |
| East Asian | | |
| With Makeup | 80% | 66.66% |
| Without Makeup | 80% | 73.3% |

*Table 3*. Accuracy of gender estimation.

To understand the cause of this bias, I created a random collection of stock images for three ethnic groups, with portraits of 30 women with makeup and 30 without makeup in each group. I added a control group of down-scaled versions of these portraits to reveal the effect of image resolution on CV biases (*Table 3*). As the comparison suggested, the gender estimation model of OpenFace was generally biased toward males, and women without makeup were more likely to be labelled as men. The reduction of resolution had the greatest impact on the recognition of Asian women. It implied that the gender estimation model of OpenFace, which was trained on iMBD's face dataset, expected more stereotypical features to recognise women and mispresented Asian women in low-resolution images.[33] As the low resolution of the films used in my project might exacerbate this computational gender bias, I did not include the data generated with this model in my research.

---

[33] 'IMDB-WIKI - 500k+ Face Images with Age and Gender Labels'
<https://data.vision.ee.ethz.ch/cvl/rrothe/imdb-wiki/> [accessed 14 June 2023].

## 4. Summary

My project has practically explored the use of computer vision techniques in the analysis of a large corpus of films. Based on the rigorous use of documents from film archives and unofficial sources, I created a visual corpus of 732 films to support my methodological experiments. While this large corpus of moving images has enabled me to incorporate more CV techniques into the research of visual culture, my practical experience in this project has also raised cautiousness in two aspects. First, the use of data provided by official film archives should be based on the awareness of the potential limitations caused by political and historical contexts. Although I combined several sources to collect the metadata for the corpus, the online database 1905.com was used as the only source for video data collection to ensure consistency. This database had limited coverage of the films produced by private institutions in Shanghai before 1953, and the relatively low resolution of the video data provided by this website has limited my experiments with the methodology. Second, because the research of large-scale visual corpus relies heavily on computer vision (CV) techniques, a deeper understanding of the potential biases caused by the application of CV would be necessary. The vision of computers is not only a substitute for the human eye but also a pixels-based system of visual representation. Therefore, critical engagements with CV techniques are a prerequisite for large-scale studies of visual culture.

# Bibliography

Arnold, Taylor, and Lauren Tilton, 'Distant Viewing: Analyzing Large Visual Corpora', *Digital Scholarship in the Humanities*, 34.Supplement_1 (2019), i3–16

Arnold, Taylor, Lauren Tilton, and Annie Berke, 'Visual Style in Two Network Era Sitcoms', *Journal of Cultural Analytics*, 4.2 (2019)

Arnold, Taylor, Lauren Tilton, and Justin Wigard, 'Automatic Identification and Classification of Portraits in a Corpus of Historical Photographs' (presented at the Computational Humanities Research Conference 2022, Antwerp, Belgium: University of Antwerp, 2022), pp. 25–35

Baltrusaitis, Tadas, Amir Zadeh, Yao Chong Lim, and Louis-Philippe Morency, 'OpenFace 2.0: Facial Behavior Analysis Toolkit' (presented at the 13th IEEE International Conference on Automatic Face & Gesture Recognition, Xi'an: IEEE, 2018), pp. 59–66

Cai, Kevin G., *Cross-Taiwan Straits Relations since 1979 Policy Adjustment and Institutional Change across the Straits* (Singapore; Hackensack, NJ: World Scientific, 2011)

China Film Archive, *中国影片大典：故事片，舞台艺术片，1949.10-1976 [Encyclopaedia of Chinese Films, 1949.10-1976]*, 1st edn (Beijing: China Film Press, 2001)

'cinematic', 2023 <https://github.com/El-Mundo/cinematic> [accessed 5 January 2025]

El-Keilany, Alina, Thomas Schmidt, and Christian Wolff, 'Distant Viewing of the Harry Potter Movies via Computer Vision' (presented at the The 6th Digital Humanities in the Nordic and Baltic Countries Conference, Uppsala, Sweden: Uppsala University, 2022), pp. 33–49

Fabbrizzi, Simone, Symeon Papadopoulos, Eirini Ntoutsi, and Ioannis Kompatsiaris, 'A Survey on Bias in Visual Datasets', *Computer Vision and Image Understanding*, 223 (2022), 103552

Firmansyah, Andrian, Tien Fabrianti Kusumasari, and Ekky Novriza Alam, 'Comparison of Face Recognition Accuracy of ArcFace, Facenet and Facenet512 Models on Deepface Framework', 2023, pp. 535–39

Heftberger, Adelheid, *Digital Humanities and Film Studies: Visualising Dziga Vertov's Work*, Quantitative Methods in the Humanities and Social Sciences (Cham: Springer International Publishing, 2018)

'IMDB-WIKI - 500k+ Face Images with Age and Gender Labels' <https://data.vision.ee.ethz.ch/cvl/rrothe/imdb-wiki/> [accessed 14 June 2023]

Jang, Ji Yoon, Sangyoon Lee, and Byungjoo Lee, 'Quantification of Gender Representation Bias in Commercial Films Based on Image Analysis', *Proceedings of the ACM on Human-Computer Interaction*, 3.CSCW (2019), 198

Jin, Fenglan, and Zhengduo Yang, 中国艺术影片编目 *[The Catalogue of Chinese Artistic Films]* (Beijing: Culture and Art Publishing House, 1981)

Jocher, Glenn, Ayush Chaurasia, and Jing Qiu, 'YOLO by Ultralytics', 2023 <https://github.com/ultralytics/ultralytics> [accessed 14 June 2023]

Jofre, Ana, Vincent Berardi, Carl Bennett, Michael Reale, and Josh Cole, 'Faces Extracted from Time Magazine 1923-2014', *Journal of Cultural Analytics*, 5.1 (2020)

Kane, Daniel, *Chinese Language: Its History and Current Usage*, 1st edn (North Clarendon, VT: Tuttle Publishing, 2006)

Lansdall-Welfare, Thomas, Saatviga Sudhahar, James Thompson, Justin Lewis, FindMyPast Newspaper Team, and Nello Cristianini, 'Content Analysis of 150 Years of British Periodicals', *Proceedings of the National Academy of Sciences*, 114.4 (2017), E457–65

Männistö, Anssi, Mert Seker, Alexandros Iosifidis, and Jenni Raitoharju, 'Automatic Image Content Extraction: Operationalizing Machine Learning in Humanistic Photographic Studies of Large Visual Archives', in *arXiv: Computer Vision and Pattern Recognition*, 2022

Marion, Donald J., *The Chinese Filmography: The 2444 Feature Films Produced by Studios in the People's Republic of China from 1949 through 1995* (Jefferson, NC; London: McFarland & Co, 1997)

Mazières, Antoine, Telmo Menezes, and Camille Roth, 'Computational Appraisal of Gender Representativeness in Popular Movies', *Humanities and Social Sciences Communications*, 8 (2021), 137

Michel, Jean-Baptiste, Yuan Kui Shen, Aviva Presser Aiden, Adrian Veres, Matthew K. Gray, The Google Books Team, and others, 'Quantitative Analysis of Culture Using Millions of Digitized Books', *Science*, 331.6014 (2011), 176–82

Morris, Travis, 'Extracting and Networking Emotions in Extremist Propaganda' (presented at the 2012 European Intelligence and Security Informatics Conference, Odense, Denmark: IEEE, 2012), pp. 53–59

Park, Seho, Kunyoung Lee, Jae-A. Lim, Hyunwoong Ko, Taehoon Kim, Jung-In Lee, and others, 'Differences in Facial Expressions between Spontaneous and Posed Smiles: Automated Method by Action Units and Three-Dimensional Facial Landmarks', *Sensors*, 20.4 (2020), 1199

Pi, Binyan, '中国"十七年"电影研究的回顾与反思 [A Review and Reflection of Research on Chinese Cinema in the "Seventeen Years Period"]', *Bohai University Journal (Social Science Volume)*, 39.2 (2017), 100–104

'Romanised Names', 2023 <https://github.com/El-Mundo/cinematic/blob/master/OCR/JAVA/romanizing/roman-non-han_names.csv> [accessed 5 January 2025]

Schroff, Florian, Dmitry Kalenichenko, and James Philbin, 'FaceNet: A Unified Embedding for Face Recognition and Clustering' (presented at the Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, Boston, MA: IEEE, 2015), pp. 815–23

Serengil, Sefik Ilkin, and Alper Ozpinar, 'An Evaluation of SQL and NoSQL Databases for Facial Recognition Pipelines' (Cambridge Open Engage, 2023) <https://doi.org/10.33774/coe-2023-18rcn>

'SimpleFaceClustering', 2023 <https://github.com/El-Mundo/SimpleFaceClustering>
[accessed 5 January 2025]

Wang, Zhuoyi, *Revolutionary Cycles in Chinese Cinema, 1951-1979*, 1st edn (New York: Palgrave Macmillan, 2014)

Zhang, Yingjin, *Chinese National Cinema*, 1st edn (New York; London: Routledge, 2004)

Zhao, Tongsheng, 电影档案概论 *[Introduction to Film Archiving]* (Beijing: China Radio, Film and Television Press, 2005)