

Metahistory for (Ro)bots: Historical Knowledge in the Artificial Intelligence Era

Meta-história para robôs (bots): o conhecimento histórico na era da inteligência artificial

Thiago Lima Nicodemo & Oldimar Cardoso

<https://orcid.org/0000-0002-1588-0683> 

<https://orcid.org/0000-0002-5614-4535> 

ABSTRACT

This text offers a theoretical reflection on the effects of the artificial intelligence and digital era on the historian's métier. It is based on a set of experiments involved in the development of a cybernetic historian, dealing with hypotheses such as (ro)bots creating historical narratives and mastering methods of both quantitative and qualitative analysis, as well as suggesting research problems. In other to do so, we present our own technology, in progress of development, and we problematize the steps to create a historian "bot". The term robot is understood as a computer program executing tasks on a largely automated basis, without any relationship with a human user. In turn, tasks are complemented by an artificial intelligence system. This emergent reality raises an urgent debate on ethical issues, such as transparency and digital ethics, and it may also be useful to problematize the future of the historical profession in the contemporary world.

KEYWORDS

Digital Humanities; Digital History; Theory of History.

RESUMO

Esse texto oferece uma reflexão teórica sobre os efeitos da inteligência artificial e do universo digital no ofício do historiador. A reflexão é baseada em um conjunto de experimentos relacionados com o desenvolvimento de um "historiador cibernético", lidando com hipóteses tais como, robôs criando narrativas históricas e dominando métodos de análise qualitativa e quantitativa. Para isso, apresentamos nossa tecnologia própria em fase de desenvolvimento, problematizando as etapas para a criação de um "robô" historiador. O termo "robô" (ou "bot") é entendido como um programa computacional que executa tarefas de forma quase inteiramente autônoma, sem qualquer relação com o usuário humano. Por sua vez, estas tarefas são complementadas por um sistema de inteligência artificial. Essa realidade emergente suscita questões urgentes sobre transparência e ética no mundo digital, e pode ser uma poderosa ferramenta para problematizar o futuro da história no mundo contemporâneo.

PALAVRAS-CHAVE

Humanidades Digitais; História Digital; Teoria da História.

Post-human narratives

Orson Krennic (Ben Mendelsohn), Director of the Advanced Weapons Research Division of the Imperial Forces enters the space shift and finds his commander, Governor Tarkin (Peter Cushing), standing back, looking through a wide glass window at the final stages of his planet-sized doomsday weapon, the “Death Star.” Still facing back, the Governor expresses his disappointment with the “security breach” on Jedha with a biting voice inflection somewhat resembling a classical horror movie’s butler. As Tarkin turns to face his interlocutor, the audience is surprised to see the same actor from the late 1970s first *Star Wars* saga movie appear on the screen forty years later, with no sign of aging. The spectator feels something is not right in his acting: a slightly robotic movement in his articulation, a somewhat rubbery texture on his face. And, besides, how could an actor that looked around 70 years old in 1977 not be dead in 2016?

Peter Cushing did, in fact, pass away in 1994. His post-mortem role in the *Star Wars* saga movie *Rogue One* was possible thanks to a high-tech computer-generated image overlapping a real actor’s performance (Guy Henry). In essence, such process is based on an overwhelming personal archive: a complete record of all of Cushing’s performances, including his roles as Frankenstein, Sherlock Homes, and Dracula in the 1940s and 1950s. Every facial expression, every voice inflection and body gesture, including a Cushing’s mid-eighties mask plaster lifecast, was used as input for the robot-avatar.

This remarkable technological achievement has raised attention in terms of ethics and legal dilemmas, such as limits and consents for a non-human after-death performance. For example, actress Carrie Fisher (Princess Leia), who appears forty years younger in *Rogue One* thanks to the same technology that enables Cushing’s acting, and passed away as *The last Jedi* (2017) was being produced, apparently had given consent for the use of her image in the follow-up episodes of the *Star Wars* franchise. Another suggestive example is the appearance of the

character Rachel (Sean Young) in *Blade Runner 2049* (2017). A similar technology allows the 58-year-old actress to appear exactly the same as 35 years earlier. The interesting spot in this case is the metafictional or intertextual element, since Rachel herself impersonates a droid seducing her romantic partner in the 1980s *Blade Runner*, now aged Deckard (Harrison Ford). Reality in this case matches fiction, because Young plays what she indeed is: a hi-tech avatar, based on her 35-year-younger self. The issue of personal archives or memory enabled by technology concerns not only sci-fi movie actors, but every ordinary man or woman who produces massive digital information through computers, digital media, photos, video recording, interactions with friends and colleagues, and even texts, all of which might be subject to memory reproduction and eventual impersonation at different levels.

One could thus think of a “biopolitics of memory,” an idea that cannot be taken for granted taking into consideration Foucault’s and Agamben’s writings on the topic. Agamben claims biopolitics as a fundamental concept that stresses the original bond between politics (sovereignty) and the “bare” life. His master metaphor in *Homo Sacer* is based on the linguistic difference in Ancient Greek between the meanings of life: *zoē*, the life proper to all living things, and *bios*, life in interaction, which could be understood as political life in its primordial form (AGAMBEN 1998a). Control over a biological body, even when stripped of its political qualities, such as the figure of the banning or the “Muselmann” (AGAMBEN 1998b, p. 155) (a nazi concentration camp refugee figure described by Primo Levi, in which violence and malnutrition leads to a state of bare life latency), sets the original source of sovereignty within modern States: control over bodies (AGAMBEN 1998b).

Thinking in Agambean terms, the biopolitics of memory implied in Cushing’s after-death performance corresponds to the widening of biopolitical control beyond bodies as well as beyond death; it neutralizes the dichotomy between *zoē*, and *bios* and enables the possibility of a post-human paradigm: the

control over a bare *bios*, a political control over minds, even in the absence of the body. The suggestive idea of a mind with no body being controlled remotely has grown into a ubiquitous futuristic topos within the sci-fi genre that can be defined as “mind uploading” or “whole brain emulation.” Variations on the same theme can be found in series and movies, such as the British production *Black Mirror*. A recent episode, “Black Museum” (2017, Season 4, Episode 6) captures the idea by presenting a collection of crime stories related to “mind uploading” – the transfer of a mind into a device or another being. In the story, the spectator realizes that some of the artifacts collected are the trapped minds of the very people involved in those crimes. In some cases, such as in *Black Mirror*’s “USS Callister” (2017, Season 4, Episode 1), the mind uploading results in a duplicate consciousness, in this case trapped into a sadistic payback role-playing game. In other cases, such as on the pilot for the whole series “Altered Carbon” (2018), man reaches immortality thanks to the shifting of someone’s mind to different bodies through a mini-disc stored in the back of the head.¹

This text presents a series of experiments dealing with an analogous idea: the possibility of a non-human writing of history, enabled by a computer program and a very detailed input or archive. In short, the historian bot would operate in a somehow similar way to Peter Cushing’s avatar, or any other sci-fi analogy mentioned so far. The development of a computer program capable of processing historical information and producing texts is not the main goal of the project hereby presented nor its possible commercial applications. In other words, we are not hoping that a historian bot will be fully functional anytime soon, but it must be seen as a hypothetical horizon.² However, in order to deal with this hypothesis, this text will problematize the concrete steps for a historian bot to be successfully functional, and, at the end of the text we will show a complete flowchart and an item, entitled “The Algorithm”, entirely dedicated to explain the technical steps of the bot. In the course of the analyses we will also refer non-systematically to some of the algorithm key steps. In

1 - Other examples of recent series are Black Mirror, Season 2, Episode 4, “White Christmas” (2014) and X-Files, Season 11, Episode 2, “This” (2018).

2 - We understand both robots and bots as programmable things that execute actions automatically. The term “bot” derived from “robot” just because of the corporeal culturally attributed characteristics of the robots. Our proposition intentionally plays with this meanings. It is also worth pointing out that the research presented in this text not only reflects about technology but is directly involved in the creation and development of new technology. There is no “software” used in these experiments: creating a bot is more complex than using a software, requires writing a complex code in a programming language, which in this case is Python.

addition, some of the bot steps can help developing effective researching tools for historical research, as shown further. As Manovich, Silveira and others assert, the digital media emerges as a transposition of traditional media and data into computer programming language (MANOVICH 2001, p. 46-47; SILVEIRA 2018, p. 106-108). Digital media is, therefore, a cultural form with a strong claim of objectivity regarding making meaning out of the world (GALLOWAY 2012, p. 54-77). Research tools in digital humanities operates according to this very same logic in transposing traditional data into digital forms associated to a rhetoric of objectivity. This text is inspired by an idea of a possible “metahistory” of the digital research tools that might be useful for the historian’s craft. Of course, it takes this idea from the well-known Hayden White’s book, “Metahistory”, a book moved by the idea of scrutinizing the discourse structures and implicit rules underlying the XIXth century European historical imagination.

In both metahistorical cases objectivity as a rhetorical form plays a fundamental role (WHITE 1973, p. 433-434). As Ramsay argues, the frame for investigation should be the “hermeneutical foundations that make such statements seem necessary” instead of “the nature and limits of computation (which is mostly a matter of methodology) and move it toward consideration of the nature of the discourse in which text analysis bids participation” (RAMSAY 2011, p. 8). In any case, the database is a “cultural form” very resisting to interpretation because refuses to project a previous order to the world of meanings (MANOVICH 2001, p. 225). This form deeply contrast with the traditional forms of history and literature understood as “narrative”, because what makes a narrative is the organization of apparently chaotic events in a plot (RICOEUR 1983, cap. 2). Therefore, as Manovich asserts, digital database and narrative are concurrent forms, “natural enemies”. In his own words, “competing for the same territory of human culture, each claims an exclusive right to make meaning out of the world” (MANOVICH 2001, p. 225). In resume, the metahistorical horizon of the XXIth century must consider the tension between

narrative and digital databases (as a cultural form) as not only valid but as a fundamental question.

Moreover, in order to substantiate or claim for full transparency from the “historian robot” in its own making as a software that may be used as an educational and learning tool. As the code created for this robot would be written as a reflection of what historical knowledge is and what a historian does, this code is itself understood as a new metahistory, or at least could help provide new grounds for future metahistorical exercises.

The input: on “hyper-archives”

The crystallization of the “mind uploading” *topos* in recent sci-fi might be understood as a symptom of significant shifting within the genre. Fiction in the last decades of the twentieth-century, such as the book *Do Androids Dream of Electric Sheep?* (1968) and the movies *Blade Runner* (1982) and *Terminator* (1984), dealt with the dystopian fear of humans being replaced by robots. Authors were constantly driven by the idea that control over society is lost to androids that eventually identify humanity’s remains as threats or flaws.³ In sharp contrast, mind uploading narratives rely on the possibility of full control over minds even in the absence of bodies or, turning again to the Agamben-inspired idea, a “biopolitics” of memory. The common element to these narratives is no longer the struggle between man and machine, but the very idea of scraping and storing unlimited personal information. Again, mankind fails and loses control to a cybernetic will, yet there is no embodiment on either side, man or machine, but rather a non-visible threat, underlying our experiences in everyday life, with searches on the web, social media, email, etc. Thus, the main question underlying the “deep” use of artificial intelligence in the new wave of robotics is: how do we define something that we cannot see but has great control over our lives, such as Facebook or Google?

3 - Although, as Hayles notes, the drama is sustained by a dialectical drive between the human element inside the non-human and vice-versa. See Hayles (1999).

The “mind uploading” *topos* shows then a displacement of the biopolitical focus from bodies to minds, a reaction to the

sense of ubiquity of control in contemporary society. Most of Agamben's work deals with the opposite possibility, which is the partial or total death of the social and political capabilities of the "animal on the outside," relegating the biological body, "the animal on the inside," to political control (AGAMBEN 1998b, p. 152) Mind uploading deals with the survival of life beyond biological restraints, the body. In his own words, "whether what survives is the human or the inhuman, the animal or the organic, it seems that life bears within itself the dream – or the nightmare – of survival" (AGAMBEN 1998b, p. 155).

What is at stake is no longer the coming of the "terminator" to annihilate humanity, but rather of an invisible algorithm or artificial intelligence that affects our lives very deeply, amplifying a sense of surveillance and lack of privacy. The "mind uploading" *topos* is nothing but a metaphor for this invisible threat, the symbol of an impossible disarticulation of the subject beyond contingency and possibility (AGAMBEN 1998b). Mind uploading *topoi* go even beyond Hayles' definition of post-human as "data made flesh" (quoting Gibson's *Neuromancer*), but in a post-biological direction – essentially, flesh made data (HAYLES 1999, p. 5-6).

Turning back to Cushing's performance as an example, one could argue that what makes his avatar plausible is the mobilization of an overwhelming repository of personal information; on another plane, closer to reality and everyday life, it can be argued that the capabilities of making significant correlations within Google or Facebook, which have contributed so much to raising the feeling that we are no longer in control of our lives, are also made possible for the very same reason. Assuming we are dealing with forms of archiving and storing information, the question raised by this statement is: should archives be re-conceptualized, considering these new social outcomes? The answer to this question can help lead the discipline of history and historians to the frontline of social science research, or at least allow for rethinking some fundamental aspects of its epistemology, since documents and

archives have always been a central foundation of nineteenth and twentieth century historical research (WIMMER 2015).

Traditionally, archives are the physical place where data is accumulated, after a process of collection, conservation, and classification. Every archival system has a “threshold,” a point at which an archive takes physical custody of records. Normally, this threshold is regulated by a “retention schedule,” a set of rules established by the archive to assess what is going to be permanently stored or disposed of (PEARCE-MORSES 2005; SCHELLEMBERG 1996). When a document becomes permanent or historically relevant, it loses its original function (which implies transformation), and that is why the retention schedule is specific to the context where the document flow occurs.⁴

Ricoeur and de Certeau consider the archive to be not only a physical place, but a “social place” as well. In Ricoeur’s words,

the multileveled architecture of the social units that constitute archives calls for an analysis of the act of placing materials in such archives, their archiving, capable of being situated in a chain of verifying operations [...] (RICOUER 2006, p. 167).

4 - For a brief history of the archive, see Giannachi (2016, p. 1-25).

There are social protocols underlying the cognitive operations implied in archiving or, expanding this argument in Foucault’s terms, “the general system of the formation and transformation of statements” (FOUCAULT 1972, p. 130).

Technological and communication processes in the contemporary world produce massive quantities of historical data and might be understood as archives in both terms: physical storage and social entity. An archive means at the same time the physical storage and its power of consignment, a set of rules and social protocols that merge into a system of signs and meanings. But the question is: what is a hyper-archive? Is there any differentiation from regular archives, considering the duality in every archive? Hiperarchives can be, as Cohen and Rosenzweig (2011) asserts, far larger, more diverse and more

inclusive than traditional archives. That is because “with new media, the content of the work and the interface are separated. It is therefore possible to create different interfaces to the same material” (MANOVICH 2001, p. 227). The archive is not only the input, but, quoting Manovich, the “center of the creative process in the computer age”. Forgetting is as constitutive of memory as disposal is of an archive. What makes a hyper-archive different from a traditional archive is precisely the loss of agency over forgetting and erasing, which ultimately results in someone or something living in a loop, not being allowed to die. We don’t even need to resort again to Cushing’s after-death performance as a metaphor; just consider the controversy over erasing information on Facebook, Google, etc. Writing about forgetting as a constitutive element of memory, Ricoeur inquires, “could a memory lacking forgetting be the ultimate phantasm, the ultimate figure of this total reflection that we have been combatting in all of the ranges of hermeneutics of the human condition?” (RICOEUR 2006, p. 413). Having in mind the case of Borges’ *Ficciones*, “Funes el memorioso” (BORGES 1988), a man incapable of forgetting anything, Ricoeur defines this question in terms of a feeling, a “presentiment” (“*Ahnung*”), “as we pass through the procession of figures that hide the horizon line.” (RICOEUR 2006, p. 413).

Getting back to Agamben’s biopolitics, the formula “to make live and to let die” is the “insignia of biopower” and it differs itself from the dynamics of the sovereign power in the old territorial State defined by Foucault, summarized by the *formula* “to make die and to let live.” The reflection on contemporary mind uploading narratives leads not “to make die or to make live, but to make survive,” still following Agamben’s definition of twentieth century biopolitics; in his own words, “the decisive activity of biopower in our time consists in the production not of life and death, but rather a mutable and virtually infinite survival.”(AGAMBEN 1998b, p. 55). The basic difference between Agamben’s biopolitics and the implications of mind uploading fictions/hyper-archives is that, in the former, *zoē* and *bios*, the inhuman and the human, are disrupted

through an emphasis on the biological body, whereas in the latter the emphasis relies on the political or social body, through a hypertrophy of data and memory. It can be argued that the hyper-archive gets even closer to the biopower's supreme ambition: "the absolute separation of the living being and the speaking being, *zoē* and *bios*, the inhuman and the human – survival" (AGAMBEN 1998b, p. 156).

Finally, as for the historian robot experiment, it can be stated that what feeds every robot is nothing but a hyper-archive, a digital documentation and/or bibliographical corpus. Thus, the basic principle of a historian robot is in fact data scraping. Our particular experiment is based on simple material scanning, followed by Optical Character Recognition (OCR) application. This procedure corresponds to the fifth step of the historian robot (flowchart box 5). The bot also depends on a careful text preparation of the sources (flowchart box 6) through data treatment by deleting duplicate pages and junkpages (such as advertising and tables of contents), merging portable document format (pdf) files (to combine many articles in just one full edition file), converting these files to .txt format, converting these files' system from Unix to DOS, deleting headers and footers from each page (as they involve repetition of the same words on many pages, which can skew the final word count and all the results), and merging the .txt files to create the corpus of each selected period. Then, the archives to be processed can be seen in the two senses already mentioned: as storage and as power of consignation, a set of rules and social protocols that merge into a system of signs and meanings.

However, such a process could theoretically be applied to any digital book or document database. That is why scanning projects should be pushed further, but bearing in mind that, in the short term, such "archives" could undergo massive robotic assessments. In other words, there is far more room for knowledge to be produced once archives become digital, as we show in the next item.

Heuristics of the new times

The historian robot idea represents a displacement of the technology originally conceived to trace consumer profiles towards production of historical knowledge. This shifting is done simply by feeding the robot with an archive to be processed. In order to explore the potentialities of these tools for our purposes we must consider the actual reach of artificial intelligence technology beyond tracing consumer profiles. For instance, displacing this original function is analogue to what the consulting company Cambridge Analytica did in 2014, by directly and indirectly collecting personal information from 50 million people through Facebook (RILEY; FRIER; BAKER 2018). Instead of tracing consumer profiles, the company used psychographic modeling techniques to generate political profiles that might have been used to target voters during the United States' 2016 presidential campaign. According to Michael Riley and others, "the firm believed those profiles were better predictors of how voters could be swayed through targeted ads than traditional data on party registration and voting patterns" (RILEY; FRIER; BAKER 2018). It is possible that the same company influenced the Brexit vote in 2016, by identifying masses of voters more susceptible to manipulation.

However, such wide-ranged technologies generally understood as "text mining" could also be a very powerful tool in scientific research if questions are asked considering other types of databases. Text mining tools are able to compute lexical patterns in frequency and distribution of words and performing tasks as grouping and categorization (JOCKERS 2013, p. 24-34). Very recent initiatives, for example, apply a knowledge-graph-based system in the probabilistic search for adequate drugs for cancer treatment (MCCUSKER et al., 2017). Recent experiments on AI conducted by Caliskan et al. (2017) at Princeton University developed a word-embedding method algorithm capable of representing each word in its interactions within a text corpus of 2.2 million unique words (out of 840 billions of tokens) and in 300 semantic dimensions (named WEAT - Word-Embedding

Association Test). The result shows not only that AI devices incorporate human-like biases, such as gender prejudice and others, but also that human prejudicial behavior implicitly conveys ingroup/outgroup identity information through language (CALISKAN et al. 2017).

Raw linguistic evidence, amplified on an unprecedented scale, confirms what we already knew from twentieth century linguistics: firstly, that meaning is defined by use; in other words, that there is a visceral correlation between meanings and speech acts (WHORF 1956), and, ultimately, “that behavior can be driven by cultural history embedded in a term’s historic use” (CALISKAN et al. 2017, p. 185). But the huge difference of scale allows for relevant progress since we can much better assess the intricate system of meanings where a word is embedded. It can be argued that this intricate network of correlations is nothing but an archive. Following Agamben’s reading of Foucault, the archive corresponds precisely to the threshold between meaning and speech, so “the archive is thus the mass of the non-semantic inscribed in every meaningful discourse as a function of its enunciation,” and furthermore, that the archive reduces the subject to a “simple function or an empty position,” (AGAMBEN 1998b, p. 145) and it “is the unsaid or sayable inscribed in everything said by virtue of being enunciated; it is the fragment of memory that is always forgotten in the act of saying ‘I.’” (AGAMBEN 1998b, p. 131).

An engine such as this one works basically by deriving artificial intelligence “by discovering patterns in existing data.” (CALISKAN et al., 2017, p. 183). This can be made for example by topic modeling, a type of statistical modeling for discovering the abstract “topics” that occur in a collection of documents. There is no predictable feature to patterns revealed in processing massive linguist evidence. It is also hard to find a graphical representation of multiple overlapping dimensions, including the modification over time and space. Generally, the idea of “network” is for semantic networks and eventually knowledge graphs (BRACHMAN 1979). Variations of this technology have

been developed since the 1950s and have spread in recent years through commercial and research applications on the web, such as Google's "knowledge graphs" since 2012 (ROUSH 2012). In the experiments discussed in this paper we have used something simpler than topic modeling, which is the counting of the most used words from a corpus under Zipf's Law (ZIPF 1949) and the methodology described by Silva and Silva (2016). This is an empirical law on mathematical statistics, which determine that the frequency of any word in an ordered list is inversely proportional to its rank in the frequency table. A word is less relevant in a corpus the more advanced is its ranking position, the majority of the words have very low frequencies and play an irrelevant role in it. The decreasing of the relevance of each word on the ordered list is often logarithmic, instead of linear, so the most used words in a corpus are completely relevant to establish its essence.

On the graph below we show an example of semantic network based on a single corpus analyses experiment: all the texts published by the journal *The Public Historian* during its thirty years of existence, from 1978 to 2017.

scrutiny of the nation-state oriented history? In any case, the way of collecting data, displaying this data and interpreting it, totally depends on the on the human eye. The machine, so far, only enhances the capacity of data processing.

Another fundamental step of the historian robot works by creating a series of semantic networks over pre-determined time frames. Technically, it starts with the lemmatization (the algorithmic process of determining the lemma of a word based on its intended meaning, e.g. by grouping together the inflected forms of a word) of the corpora, which differentiates nouns and verbs written with the same words and divides compound words (flowchart box 7). After that, the words of each corpus are counted and ranked in the file `words.csv`, the nodes formed by these words are identified by one identity number in the file `nodes.csv`, and these identity numbers are used to establish the edges among these words. After this, we proceed to a manual input of stopwords, which are the non-relevant words to the research (such as *the, of, and, be, to, etc.*) (flowchart box 8), the words used in the title of the source (such as the words *public* and *historian* in the case of the journal *The Public Historian*) and in the title of the field research related to the source (like the word *history* in the case of the journal *The Public Historian*). The bot can be loaded with a generic list of stopwords and skip this manual step by automatically filtering the ranking of the most used words using its default list of stopwords, but manual input provides better quality until a specific artificial intelligence (like the Application Programming Interface `spacy.io`) is developed to define the stopwords of each corpus. Then the robot assesses if there still are stopwords among the words on the file `words.csv` (flowchart box 9), which enables the sorting of the most used words in different periods by merging all the semantic networks of each period in one temporal network (flowchart box 10),⁵ and the consecutively summing junction which uses the equation $R = (k \times w) + \ell + M + \mu + \tilde{x} + m + s + \Delta + v + a + \sigma$ to merge the rankings of the most used words in the corpora of different periods (flowchart box 11). For example,

5 - To the concept of temporal network, see: Peixoto e Rosvall (2017) and Li et al. (2017).

such equation could be used to merge the rankings of the most used words in the corpora of five decades, from the 1970s to the 2010s. In general terms, the results could be effectively used to assess popular themes or trends in the historiography of the “Public History” fields. Similar methods could be applied to historical sources or any kind.

Table 1 – The most used words in all editions of the journal *The Public Historian*

	1978-2017	R	1970s	1980s	1990s	2000s	2010s
1 st	American	American	historical	historical	American	American	museum
2 nd	historical	national	program	state	historical	national	American
3 rd	work	work	work	work	national	historical	national
4 th	museum	state	university	American	work	work	work
5 th	national	historical	state	national	state	state	state
6 th	state	museum	research	program	museum	university	park
7 th	university	university	study	university	book	site	historical
8 th	park	park	policy	Study	university	exhibit	site
9 th	book	site	project	Policy	study	park	university
10 th	time	time	student	Book	time	Time	community
11 th	site	war	department	Research	war	People	project
12 th	program	book	national	Record	site	Book	war
13 th	war	commu- nity	people	Social	park	War	time
14 th	people	people	American	Time	research	Past	people
15 th	research	exhibit	preservation	government	program	commu- nity	city
16 th	community	past	government	Service	people	Library	past
17 th	past	city	city	Local	review	Place	visitor
18 th	city	project	time	City	exhibit	City	place
19 th	preservation	study	business	Society	past	Visitor	exhibit
20 th	local	program	community	community	preserva- tion	Press	story

Some of the already mentioned questions could be complexified and proposed in different angles with the help of the time frames. The point, for the sake of this text argument,

is to imagine the possibility of analyzing almost infinite data, including documents and books, but also human interactions, economy, images in a scale the human eye cannot simply perceive. Thus, the advantage of the use of a robotic-made temporal network in a historical interpretation is comparable to the use of a microscope instead of the naked eye in natural sciences. The robot cannot interpret the sources better than a human historian, but this human historian might do a better work with the help of the robot. So the historian robot is more an exoskeleton than an automaton. It will not replace historians, but perfect their work. Also, this bot can help humans to introduce reproducibility in the humanities. If different historians use the same temporal network as the basis for an analysis, it is easier to establish distortions and biases.

On Digital Ethics and Learning Tools

In engineering, a black box is a system only accessible in terms of its input, output, and transfer functions, without any knowledge of its internal workings. The main current experiments with artificial intelligence or machine learning using neural networks exclusively, especially generative adversarial networks (GANs), tend to work in this way. Such assertion leads us to formulate a base law of “humanistic” robotics, inspired by Isaac Asimov’s “laws of robotics” (ASIMOV 1950, p. 40). The basic principle is that a historian robot must never be just a black box (Law number 1) in order to work with transparency. A historian robot must openly describe every step it took (Law number 2) and, for the sake of the present research, that is exactly what is done in the appendix (on the algorithm and the Metahistory Flowchart). Finally, to align the first and the second assertion, a historian bot must be able to be run on a personal computer, which makes it accessible to anyone (Law number 3). This basic set of rules may allow robots to be a self-developmental and educational tool.

Neural networks are created to relate data for which there still are no equations, to solve problems for which only the

answer is sufficient and the problem solving process is irrelevant. An example of machine learning based exclusively on neural networks are the walking bots developed by Boston Dynamics (RAIBERT et al. 2008, p. 10822-10825). They learn how to walk without any algorithmic instruction on what to do with their body, legs, or knees. They are just ordered to walk forward and have to learn by themselves how to do this. They fall for generations (and the learning of the previous generation is even transplanted to the next) until they understand how to use their body's resources to move under the effect of Earth's gravity. How they learn to walk is not important to computer engineers, as long as they learn to walk satisfactorily. As they learn by trial and error, their movements are more natural than in former robots taught to walk by lines of code describing precisely each movement.

Although how they learn is not a problem in many cases for robots, this is definitely a problem for humans. Ignoring this fact, many adaptive platforms developed with machine learning for educational purposes work as black boxes (BRUSILOVSKY; PEYLO 2003). The AI system does not care for the reasons and grounds for learning; it just recognizes in a binary way the effectiveness of the process. Computer engineers designing AI systems for education might get better results working with education experts because it would allow for a better understanding of how and why students learn better. The problem is simple: an adaptive platform which is a complete black box, which does not know why students learn better in the way they are taught, is not created to help teachers, but to replace them. Only experts can understand why it is a problem to totally replace teachers (or historians, in our case) with an artificial intelligence; computer engineers cannot. If we historians are out of this research, the writer bots and specifically the historian bots will be developed in the same way the educational adaptive platforms, without us and to replace us. As Annette Vee asserts, "treating coding literacy as a real thing allows us to anticipate this time and prepare for it with better and more inclusive educational approaches" (VEE 2017, p. 760).

Fighting the full black-box logic on historian bots is not (at least not only) a case of historians or educational experts corporatism. This is a political combat related to machine bias and to the replacement of moral (human) authority by (bot) mathematical authority. We understand that the moral authority eventually implicated in historical writing and scholarship should also be questioned. Sometimes historians' work is also a black box so the reader cannot understand exactly how certain inputs gave rise to their outputs, what are the sources, how theory and methods led to heuristics (how the sources were analyzed) and how this analysis implied the narrative (LATOUR 1999).

The issue of neutrality and objectivity is one of the fundamentals of historical scholarship. Many historical manuals beginning with Droysen's *Historik* (1854), claim that the "critique" of the sources is a fundamental step to avoid relying on the authority of texts by tradition (#33); the "chaos" of "simultaneous opinions, news, rumors" (#34); this is only the superficial origin of the historical sources. The historian must actively access biases in historical documents that make them part of their own time and space, and by doing so, as Ernst Bernheim's *Lehrbuch der historischen Methode*, by 1900, producing "self-distanciation", recognizing, as Herman Paul (2011) asserts "otherness of the past".

Massive quantitative data appears in its chaotic organization in a first regard as "independent of interpretation", nonetheless as Moretti asserts at the same time "they often demand an interpretation that transcends the quantitative realm" and, "most radically", "we see them *falsify* existing theoretical explanations" (MORETTI 2005, p. 30). This complex layers points to the underlying "assumptions about information, texts and people are "embedded in the software programs we compose" and that is why "the scrutiny of computational procedures can help us to understand the affordances and actions of the various programs on which we now depend" (VEE 2017, p. 760).

So taking into consideration the history of historical scholarship itself, we acknowledge that new hopes of

transparency in AI are directly related to the urge of documental and algorithmic critique, with the need of qualitatively situated sources and their own methods historically and socially. Applying AI to historical learning could then lead to multiple “bias catcher” robots using the available knowledge of the concepts of “eurocentrism”, following Chakrabarty (2000) e Young et al (2004) definitions, for instance.

These experiments could not only be a powerful learning tool, but also help enable new professional activities for historians based on what we have been doing at least in the last 200 years in terms of historical theory: discussing production of knowledge through analysis of sources, with particular attention to the historical biases of social groups in time frames. As Greenwald (2017) argues, technology such as WEAT could be used as a tool to “diagnose” biases in any type of media, or to associate different biases to certain social groups.⁶

The Algorithm

6 - See also Noble (2018).

In order to substantiate the laws defined in the previous section and to present an example of robotic-metahistorical reflection, we will describe the flowchart of the historian bot developed by the company run by one of this paper’s authors. The flowchart at the end summarizes all the necessary systems to perform from the treatment of the sources to the writing of the historical narrative. Each paragraph below is related to one of the boxes used in the flowchart, numbered from 1 to 19.

The starting step of the historianbot.org (flowchart box 1) is to collect the sources, by scanning printed books with some human help or by scraping data alone on the internet. The easiest, cheapest, and more effective way to scan a book, with better results on optical character recognition (OCR), is to shear its spines and scan it to a portable document format (pdf) file as single sheets. We can rebind the book after this process and make it brand new without any loss or throw it away as recycled paper. There are some cheap scanners that

could get the full text of a book by optical character recognition (OCR) just after scanning the printed pages. All historians can have such a scanner at home, without the need for expensive scanners which could only be bought by institutions. This is very important for free research. The only human work to scan a book is to shear its spines and insert at most 100 pages at a time into the scanner. The scanner can automatically collect these 100 pages one by one and create a portable document format (pdf) file with the full text as metadata supplied by its own optical character recognition (OCR) software for many languages. It may seem counterintuitive, but if the sources are already digital, the work can be harder than with printed sources. The first problem with obtaining digital sources is to scrape them from the internet. Many journals, books, or documents are not easily accessible. They can be read by humans, page by page, but it is commonly difficult to download the whole data, which is necessary for historian bots. They need all files on the drive to manage them – it is not possible to just read pages on a browser like a human. As many scholarly platforms have protection against bots, which is strange and symptomatic, data scraping requires the use of some application programming interfaces (APIs) to bypass these protections on the platforms where the sources are stored. The most common example of protection subject to bypassing by using APIs are the Completely Automated Public Turing test to tell Computers and Humans Apart (CAPTCHA). Some paid APIs are able to convert to text the image file of the CAPTCHA and to write this text in the expected field, simulating human action and enabling the download by the historian bot. For example, the full text of the 40 years of the journal *The Public Historian* was downloaded by *historianbot.org* in three hours; a human will need at least a week doing just this for many hours a day to complete the same task. Moreover, while a human bored with this task would probably leave some files behind, bots do not. After bypassing the protection, the second problem of the historian bot is to deal with digital files with lousy optical character recognition (OCR) because they were made a long time ago, when this technology was first out. So *historianbot*.

org is able to delete the old OCR and to generate it again with better technology. The margin of error of an old OCR, as we found on the journal *The Public Historian* from 1978 to 2000, is more than 30% (which completely undermines the analysis work of the bot). However, historianbot.org can reduce it to less than 1% by deleting and redoing the optical character recognition (OCR). In addition to text, the bot can also sort and rank pictures and videos if this is relevant to the research.

The second step of historianbot.org (flowchart box 2) is to calculate the margin of error of the optical character recognition (OCR) used to digitize the sources. The historian bot uses a spell-checking tool to know the margin of error of the digitalization of the sources by counting how many words are detected as wrong by the spell-checking tool and comparing this quantity of words with the amount of words in the whole text. If the margin of error is less than 1%, the data is sent to flowchart box 6, "Preparation: Data treatment." If the margin of error is more than 1%, the data is sent to flowchart box 3, "Or."

The third step of historianbot.org (flowchart box 3) is an or function which separates scanned sources, sent to flowchart box 4, "Manual operation: Redo scanning," from scraped sources, sent to flowchart box 5, "Predefined sources: Redo OCR."

The fourth step (flowchart box 4) is the manual operation to redo the scanning of printed sources with a margin of error superior to 1%. Historianbot.org cannot do anything if the scanning of a printed source is badly done and this is the only manual operation of this flowchart that cannot be replaced by an automatic one.

The fifth step (flowchart box 5) is the predefined process to redo the optical character recognition (OCR) if a portable document format (pdf) file source presents a margin of error greater than 1%. In this case, historianbot.org can automatically correct the problem by deleting the old OCR and by making a new one.

The sixth step (flowchart box 6) is the preparation of the

sources through data treatment by deleting duplicate pages and junkpages (such as advertising and tables of contents), merging portable document format (pdf) files (to combine many articles in just one full edition file), converting these files to .txt format, converting these files' system from Unix to DOS, by deleting headers and footers from each page (the repetition of the same words on many pages can skew the final word count and all the results), and finally merging .txt files to create the corpora of each selected period.

The seventh step (flowchart box 7) is the predefined process to create one semantic network for each period of time. This starts with the lemmatization of the corpora, which involves differentiating nouns and verbs written with the same words and dividing compound words. After that, the words of each corpus are counted and ranked in the file `words.csv`, the nodes formed by these words are identified by one identity number in the file `nodes.csv`, and these identity numbers are used to establish the edges among words.

The eighth step (flowchart box 8) is the manual input of the stopwords, which are the non-relevant words to the research (such as *the*, *of*, *and*, *be*, *to*, etc.), the words used in the title of the source (such as *public* and *historian* in the case of the journal *The Public Historian*) and in the title of the field research related to the source (such as the word *history* in the case of the journal *The Public Historian*). The bot can be loaded with a generic list of stopwords and skip this manual step by automatically filtering the ranking of the most used words with its default list of stopwords, but manual input provides better quality until a specific artificial intelligence is developed to define what are the stopwords of each corpus.

The ninth step (flowchart box 9) is to assess if there are still stopwords among the words in the file `words.csv`. If so, the list of words is sent back in a loop to flowchart box 7, "Predefined process: Create semantic networks;" if not, it is sent to flowchart box 10, Sort: Merge temporal network."

The tenth step (flowchart box 10) is the sorting of the most used words in different periods by merging all semantic networks of each period into one temporal network.

The eleventh step (flowchart box 11) is a summing junction which uses the equation $R = (k \times w) + \ell + M + \mu + \tilde{x} + m + s + \Delta + v + a + \sigma$ to merge the rankings of the most used words in the corpora of different periods. For example, if the position of the same word in the rankings of the 20 most used words in the corpora during five periods is 14th, 4th, 1st, 1st, and 2nd, the relative numbers to identify these positions are inverted to 7, 17, 20, 20, and 19. So the variables will assume the following values:

R = position of the word in the temporal ranking = 300.02
= 1st;

k = number of variables except k and $w = 10$;

w = weighted average =
$$\frac{\sum_{i=1}^n i \cdot x_i}{\sum_{i=1}^n i} = (7 \times 1 + 17 \times 2 + 20 \times 3 + 20 \times 4 + 19 \times 5) / (1 + 2 + 3 + 4 + 5) = 18.4$$
;

$$\sum_{i=1}^n i \leftrightarrow \frac{n(n+1)}{2} = 1 + 2 + 3 + 4 + 5$$
;

ℓ = last position = 19;

M = maximum position = 20;

μ = population mean =
$$\frac{1}{n} \sum_{i=1}^n x_i = (7 + 17 + 20 + 20 + 19) / 5 = 16.6$$
;

\tilde{x} = median = ordered positions: (7, 17, 19, 20, 20) = 19;

m = minimum position = 7;

s = swing = if m appears before M , then $s = M - m$,
otherwise, $s = m - M = 20 - 7 = 13$;

Δ = delta = l – first position = $19 - 7 = 12$;

v = speed = $\Delta / (n - 1) = 12 / 4 = 3$;

a = acceleration = $(2 \times v) / (n - 1) = (2 \times \Delta) / (n - 1)^2 = 6 / 4$
 $= 24 / 16 = 1.5$;

$$\sigma = \text{population standard deviation} = \sqrt{\frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2} = \sqrt{\frac{1}{5} [(7 - 16.6)^2 + (17 - 16.6)^2 + (20 - 16.6)^2 + (20 - 16.6)^2 + (19 - 16.6)^2]} = 4.92.$$

The twelfth step of the historian bot (flowchart box 12) is the predefined process of proposing historical problems. The organization historianbot.org has worked so far with five ordinary directive algorithms which analyze five parameters of the temporal network to write questions in English proposing problems to the sources. These parameters are:

1. the expressive rise of a word in the rankings of the most used words in the corpora through certain periods;
2. the expressive fall of a word in the same context;
3. the stability of (a) word(s) in the initial positions of the rankings of the most used words during all the periods;
4. the sudden rise of (a) word(s) highly ranked in the final periods without appearing in the initial periods;
5. the sudden appearance of (a) word(s) only in the final periods.

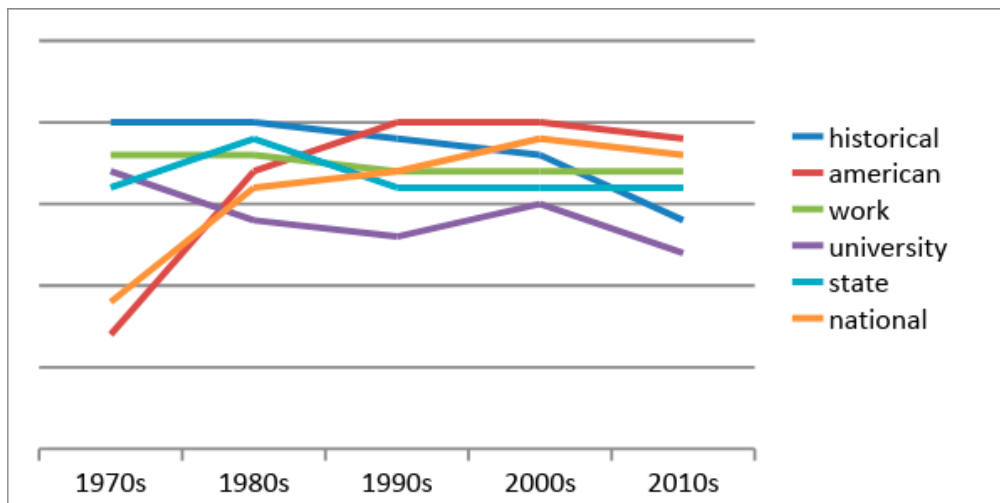
For example, the bot can formulate the following questions

to the data on Table 1 above:

1. Expressive rise: How do we explain the rise of the words *American* and *national* in the rankings of the most used words?
2. Expressive fall: None.
3. Stability: Why are the words *historical*, *work*, *university*, and *state* stable among the first half of the most used words?
4. Sudden rise: How do we explain the sudden rise of the word “museum”?
5. Sudden appearance: How to explain the sudden appearance of the word(s) *park*, *site*, *war*, *past*, and *exhibit*?

As an example of the use of this equation to merge the rankings of the most used words in the corpora of five periods, we can see on the table below the ranking of the most used words in the single corpus of all editions of the journal *The Public Historian* from 1978 to 2017 (column **1978-2017**) compared to the temporal ranking organized by this equation (column **R**) and to the five rankings of the most used words in the same journal organized by decade (columns **1970s** to **2010s**).

Graphic II – Representation of questions 1 and 3



The thirteenth step of the historian robot (flowchart box 13) is a process of writing a text to answer the questions proposed in step 12, “Predefined process: Propose historical problems” under an ordinary directive algorithm. This writing algorithm can write a text word by word using four main parameters: 1. the edges among words in quotes in the questions proposed in step 12, “Predefined process: Propose historical problems”; 2. the semantic network of the full sources; 3. the semantic network of the literature; 4. the semantic network of the author’s complete works. With these four parameters, the bot is able to establish the probability of the chain of words in a text answering each question proposed on step 12, “Predefined process: Propose historical problems.” This text is for sure still an imperfect creation, worse in style than a human research report, and demands a hard human edition.

The last six steps of the historian bot (flowchart box 14 to 19) are yet in development and are related to the creation of the historical narrative itself. A Human edition (flowchart box 14) is necessary after the Writing algorithm (flowchart box 13) to validate the narrative created by the bot based on the semantic networks. A Narrative assessment (flowchart box 15) after this Human edition decides if the narrative is ready. If it is not ready, it goes to a Neural network (flowchart box 16) and

it comes back to the Writing algorithm (flowchart box 13) to be improved. In this case, all the change decisions of the neural network are registered in a public database to provide Algorithm transparency (flowchart box 17). If the narrative is ready, it goes to the terminator as Historical narrative (flowchart box 18) and it is stored in a database (flowchart box 19) to be used in the future as part of the author's complete works to contribute to the definition of his/her text style.

Will Robots Replace Historians? Some final remarks

Just like the case of Peter Cushing's post-mortem performance, a hypothetical virtual historian could be brought back to life based on his personal hyper-archive as a source (not only personal papers, but every writing), through a mere historian's avatar. Cushing shaped his performance on his individual skills, historical circumstances, and interactions with the director and other actors. In other words, there was a unique artistic quality implied in his craft, which is lost when he becomes a "robot." The same analogy applies to a hypothetical historical robot, since it will be based on emulation and repetition of patterns. A research "methodology" (or school of thought) emulator is also very possible in the near future not only for history, but for any humanities field in general. There are already several ongoing experiments in composing music, such as the *Iamus*, at the University of Malaga⁷ (DIAZ-JEREZ 2011), the *Aiva* (Artificial Intelligence Virtual Artist),⁸ or the case of Pindar Van Arman's *Cloudpainter* robot, among many others.⁹

7 - See <http://melomics.com>.

8 - See <http://www.aiva.ai>

9 - See <http://www.cloudpainter.com>

In general, textbooks are written based on a summary of previously developed historical scholarship. Though there is room for innovation, it is oriented towards new methods of learning; in other words, on how information is displayed. For this reason, a textbook robot is very likely to be available in the short term. At least as far as the Brazilian context is concerned, textbooks are written by teams comprised of several specialists. The jobs of authors or content producers may be in jeopardy, and there will be space for a general content "curator."

Incarnating a historical method or a certain historian's style could be a very effective learning tool to help students formulate problems and enable new possibilities of working with historical sources. Moreover, the "historian robot" itself could be a heuristic tool to learn history in the digital era. Its data processing power could also be established to test old historical hypotheses and affirmations, as well as to improve existing methods in quantitative and qualitative assessments. In general terms, professional historians are losing space as agents in the production and circulation of historical knowledge in contemporary society. Digital humanities should therefore urgently be included as a discipline in historical training – however, it is worth discussing its specific conditions and above all the question: should we be turned into computer programmers?

Traditional historical training can be very useful in the Artificial Intelligence reality. We have suggested some paths, but it must also be acknowledged that experiments should be multiplied so we could understand much better these potential professional activities coming up in the near future, such as source critique robots, or "stereotype catchers." False information, diversion, "fake news" are massively replicated in social media on an increasing scale. The case of AI influencing Brexit and the US 2016 elections was particularly symbolic not only because the mass of information collected through social media enabled the tracing of behavior profiles, but particularly because, mostly, these algorithms led to locating people more susceptible to "fake news" so their opinions could be more easily manipulated.¹⁰ In other words, source critique, enabled on an unprecedented scale by AI and associated with an idea of transparency, could be a powerful tool to save whatever is left from democracy in the near future. That is one of the reasons a robot such as the historian robot must be able to run on a personal computer, must be accessed from poor or underdeveloped countries as well as open to the general public. Moreover, a historian bot must describe openly every step it took also because everyone can be able to contest its conclusions. Not only must the code used to interpret the sources be open

¹⁰ - See Hersh (2015).

to everyone, but also the full sources and bibliography used by the bot (CARDOSO 2012). This claim for transparency is strictly connected with two fundamental elements, one is the openness of the sources and codes (that the “nonprofit mission of online historical archives generally produces even higher rates of honesty” (ROSENZWEIG 2011, p. 145), and secondly with a postcolonial or peripheral horizon which fights for more equality not only for the access of digital resources but also for agency in the creation and reproduction of this same resources. Initiatives such as the “Mapping Digital Humanities in India” have shown that decentralization and empowerment in the practice of humanities brings several significant changes “particularly with respect to traditional methods of pedagogy and scholarship” (SNEHA 2016, p. 3-4).

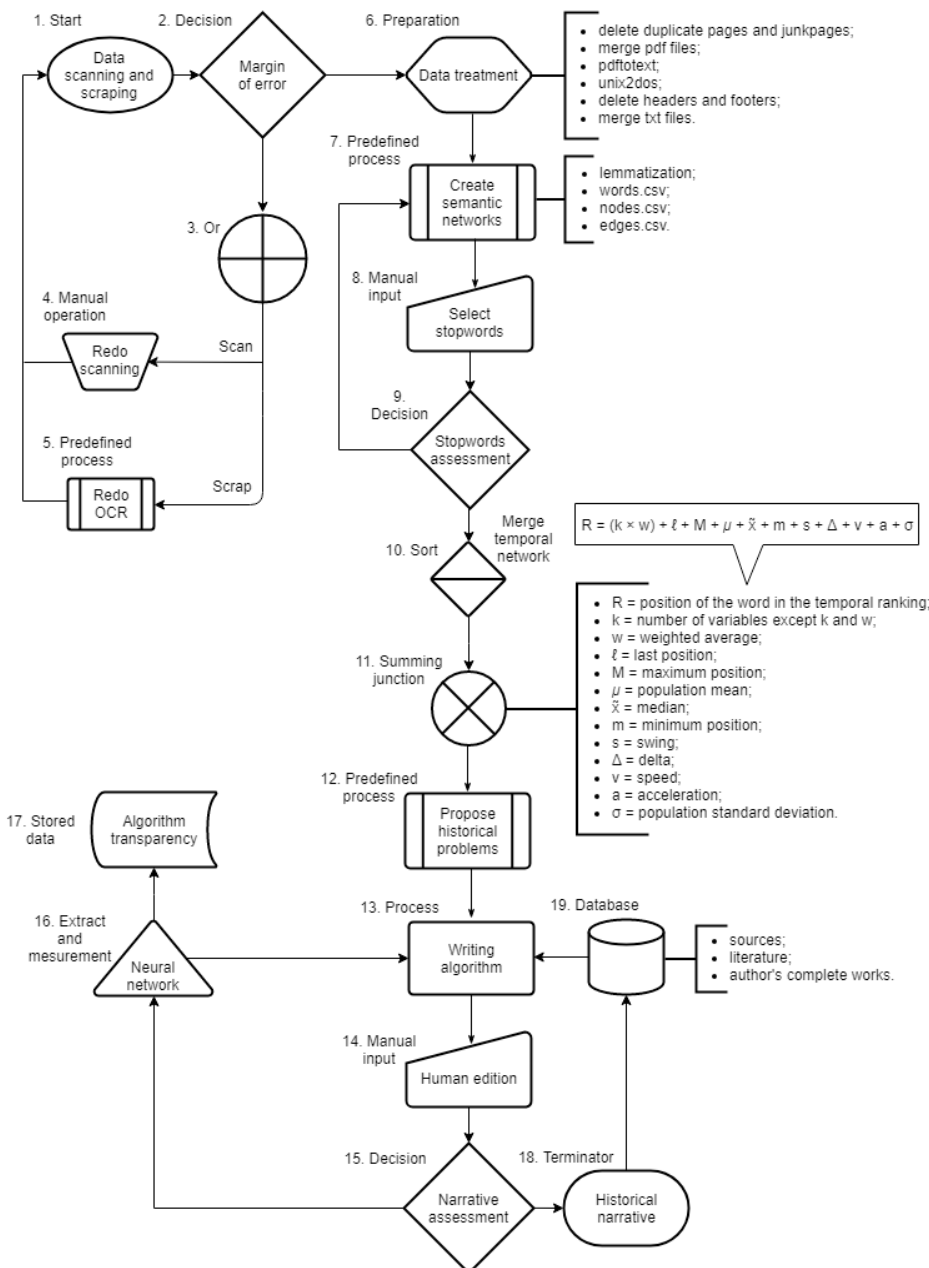
Finally, this research has also shown that digital knowledge, including libraries and historical archives, is the substantial input for post-human production of knowledge. We must know and discuss more about possible AI’s applications to digital libraries and archives. Non-human massive robotic assessments should be included in the agenda of every digitalization project. Moreover, the private monopoly of digital hyper-archives could seriously jeopardize the development of independent science (and historical knowledge) in the short or medium term. That issue raises the importance of the creation of national and global public libraries, such as the point raised in several occasions by the North-American historian Robert Darnton (2010).

Since the 1950s, cybernetics represented a threat to jobs, human ethics, and intelligence. The fact is that robots did not create their own civilization and tried to exterminate humanity as we have seen so many times in science-fiction. The real problem concerning robots – as we have learned from the WEAT (Word-Embedding Association Test) – is that robots effectively learn from humans, even unconscious prejudices and biases. Moreover, AI definitively leads everything to an unprecedented scale, including human issues such as inequality and wealth concentration, monopolies of all sorts including

knowledge, vigilance, arms races and, above all, stupidity. But it also enables some new possibilities in which historical training can still definitely contribute.

Graphic III

Metahistory as flowchart



REFERENCES

AGAMBEN, Giorgio. **Homo Sacer**: Sovereign Power and Bare Life. Stanford, Calif: Stanford University Press, 1998a.

_____. **Remnants of Auschwitz**: the witness and the archive New York: Zone Books, 1998b.

ASIMOV, Isaac. **Robot**. New York: Doubleday, 1950.

BORGES, Jorge Luis. Funes el memorioso. *In*: _____. **El Aleph**. El informe de Brodie. Caracas: Biblioteca de Ayacucho, 1988.

BRACHMAN, Ronald J. On the epistemological status of semantic networks. *In*: FINDLER, Nicholas (ed.). **Associative Networks**: Representation and Use of Knowledge by Computers. Cambridge: Academic Press, 1979.

BRUSILOVSKY, Peter; PEYLO, Christoph. Adaptive and Intelligent Web-based Educational Systems. **International Journal of Artificial Intelligence in Education (IJAIED)**, n. 13, 2003.

CALISKAN, Aylin et al. Semantics derived automatically from language corpora contain human-like biases. **Science**, n. 356, p. 183-186, 2017. Available from: <http://science.sciencemag.org/>. Accessed: Mar. 13th, 19.

CARDOSO, Oldimar. Cultura histórica e responsabilização científica. *In*: BODEMER, Klaus (ed.). **Cultura, sociedad y política en América Latina. Aportes para un debate interdisciplinario**. Madrid/Frankfurt a. M.: Iberoamericana/Vervuert, 2012.

CHAKRABARTY, Dipesh. **Provincializing Europe**, Princeton: Princeton University Press, 2000.

COHEN, Daniel J.; ROSENZWEIG, Roy. Collecting History Online. *In*: ROSENZWEIG, Roy. **Clio Wired: The Future of the Past in the Digital Age**. New York: Columbia University Press, 2011.

DARNTON, Robert, Can We Create a National Digital Library? **The New York Review of Books**, 2010. Available from: <http://www.nybooks.com/articles/2010/10/28/can-we-create-national-digital-library/>. Access: Mar. 13th, 19.

DIAZ-JEREZ, Gustavo. Composing with Melomics: delving into the computational world for musical inspiration. **Leonardo Music Journal**, n. 21, 2011.

FOUCAULT, Michel. La vie des hommes infâmes. *In*: _____. **Dits et Écrits**, Tome III, Texte 198. Paris: Gallimard, 1972.

GALLOWAY, Alexander R. **The Interface Effect**. Cambridge: Polity Press, 2012.

GIANNACHI, Gabriella. **Archive Everything**. Mapping the Everyday. Mit Press, 2016.

GREENWALD Anthony Galt. An AI stereotype catcher. **Science**, v. 356, issue 6334, apr. 2017.

HAYLES, Katherine. **How we Became Post-Human**. Virtual bodies in Cybernetics, Literature and Informatics. Chicago: The University of Chicago Press, 1999.

HERSH, Eitan D. **Hacking the Electorate**. How Campaigns Perceive Voters. Cambridge University Press, 2015.

LI, A. et al. The fundamental advantages of temporal networks. **ScienceMag**, nov., 2017.

JOCKERS, Matthew Lee. **Macroanalysis: Digital Methods and Literary History**. University of Illinois, 2013.

KELLEY, Robert. Public History: its origins, nature and prospects. **The Public Historian**, v. 1, n. 1, p. 16-28, 1978.

LATOUR, Bruno. **Pandora's Hope: Essays on the Reality of Science Studies**. Cambridge: Harvard University Press, 1999.

LIDDINGTON, Jill. What Is Public History? Publics and Their Pasts, Meanings and Practices. **Oral History**, v. 30, n. 1, Theme Issue "Women's Narratives of Resistance", 2002.

MANOVICH, Lev. **The Language of New Media**. Cambridge, Mass.: The MIT Press, 2001.

MCCUSKER, James et al., Finding melanoma drugs through a probabilistic knowledge graph. **PeerJ Comput. Sci.**, n. 4, 2017.

MORETTI, Franco. **Graphs, Maps, Trees: Abstract Models for a Literary History** Verso, 2005.

NOBLE, Safiya Umoja. **Algorithms of Oppression. How Search Engines Reinforce Racism**. New York, NY: NYU Press, 2018.

PAUL, Herman. Distance and Self-Distanciation: Intellectual Virtue and Historical Method Around 1900. **History and Theory**, v. 50, n. 107, 2011.

PEARCE-MOSES, Richard. **A glossary of archival and records terminology**. Chicago: Society of American Archivists, 2005.

PEIXOTO, Tiago; ROSVALL, Martin. Modelling sequences and temporal networks with dynamic community structures. **Nature Communications**, n. 582, 2017.

RAIBERT, Marc et al. BigDog, the Rough-Terrain Quadruped Robot. **Proceedings of the 17th World Congress The International Federation of Automatic Control**, Seoul, Korea, July 6-11, 2008.

RAMSAY, Stephen. **Reading machines:** toward an algorithmic criticism. University of Illinois Press, 2011.

RICOEUR, Paul. **Temps et Récit.** Paris: Le Seuil, 1983.

_____. **History, Memory, Forgetting.** University of Chicago Press, 2006.

RILEY, Michael; FRIER, Sarah; BAKER, Stephanie. Understanding the Facebook-Cambridge Analytica Story: QuickTake. **The Washington Post**, March 22th, 2018. Available from: <https://wapo.st/2HO1sUJ>. Access: Mar. 13th, 19.

ROSENZWEIG, Roy. **Clio Wired:** The Future of the Past in the Digital Age. New York: Columbia University Press, 2011.

ROUSH, Wade. Google Gets A Second Brain, Changing Everything About Search. **Xconomy**, December 12th, 2012. Available from: <https://bit.ly/2Ia7e2g>. Access: Mar. 13th, 19.

SCHELLEMBERG, Theodore. **Modern archives:** principles and techniques. Chicago: Society of American Archivists, 1996.

SNEHA, P.P. **Mapping Digital Humanities in India.** CIS Papers, India: The Centre for Internet and Society, 2016.

SILVA, Edson Armando; SILVA, Joseli Maria. Ofício, Engenho e Arte: Inspiração e Técnica na Análise de Dados Qualitativos. **Revista Latino-Americana de Geografia e Gênero**, v. 7, n. 1, 2016.

SILVEIRA, Pedro Telles da. **História, técnica e novas mídias:** Crítica da razão histórica digital. Tese em História Social. Programa de Pós- Graduação em História do Instituto de Filosofia e Ciências Humanas da Universidade Federal do Rio Grande do Sul, 2018.

VEE, Annette. **Coding Literacy:** How Computer Programming Is Changing Writing. Cambridge, Mass.: The MIT Press, 2017. *E-book*.

WHORF, Benjamin Lee. **Language, thought, and reality;** selected writings. Cambridge: MIT Press, 1956.

WIMMER, Mario. The Present as Future Past: Anonymous History of Historical Times. **Storia della Storiografia**, n. 68, 2015.

WHITE, Hayden. **Metahistory**. Baltimore: The Johns Hopkins University Press, 1973.

YOUNG, Robert. **White mythologies:** writing history and the west. London: Routledge, 2004.

ZIPF, G. K. **Human Behavior and the Principle of Least Effort**. Cambridge: Addison-Wesley, 1949.

AGRADECIMENTOS E INFORMAÇÕES

Thiago Lima Nicodemo

tnicodem@unicamp.br
Professor de Teoria da História
Universidade Estadual de Campinas - Brasil

Oldimar Cardoso

oldimar@gmail.com
Pesquisador no Laboratório em Rede de Humanidades Digitais
Instituto Brasileiro de Informação em Ciência e Tecnologia
Brasil

Acknowledgments: Cynthia L. Z. DeRoma, Edson Armando Silva, Marcela Guimarães Silva, Marco Costa, Pedro Telles da Silveira, Ruhama Sabião, CAPES and Alexander von Humboldt Stiftung.

RECEIVED IN: 11/JAN./2019 | APPROVED IN: 17/MAR./2019

ERRATA

Corrigendum published on April 28, 2019.

1. Both authors contributed equally to this work.
2. The property of the algorithm used to obtain the data of this article belongs to Oldimar Cardoso.