The Bitter Aloe Project: The Application of Advanced Machine Learning to the TRC Archive

The purpose of the Bitter Aloe Project is to use machine learning to bring new legibility to the TRC archive and provide new optics for inquiry into apartheid era human rights abuses. Over the past four years we have developed custom trained machine learning models to perform various forms of natural language processing on records produced by the TRC. Our named entity recognition models extracted structured data from various materials that serve as the basis for three research tools, our GIS map of human rights violations, a network graph of co-occurrences of names of individuals and organizations, and a text analysis tool based around sentence embeddings. These tools are publicly available via our website at <u>www.bitteraloeproject.com</u>. This project was supported with generous funding from the Commonwealth Institute for Black Studies (CIBS), the History Department and the Dean's Office at the University of Kentucky as well as the Harry Frank Guggenheim Foundation.

Like many truth commissions before and since, the TRC's information management approach by and largely followed what human rights practitioners call the 'who-did-what-to-whom' approach to collecting information about HRVs. This approach centers inquiry around dividing experience into a set of acts that individually serve as the most elemental datum in a data collection, emphasizing individual physical and mental harm over collective actions, and focusing the investigatory gaze on identifying individual perpetrators. The TRC, following the Promotional of National Unity and Reconciliation Act (PNURA), interpreted the scope of its mission as creating 'as complete a picture as possible' of gross human rights violations committed in South Africa between 1960 and 1994 (later extended to 1995). This picture came in the form of public testimonies, which became the primary way scholars, the press and the public engaged with the TRC's work, and less public 'statement-taking' which comprised the bulk of raw data collection, which followed the 'who-did-what-to-whom' approach.

The impact of the 'who did what to whom approach' remains a topic of debate in the literature on the TRC and in human rights studies in general. Although its proponents argue that it offers the most direct route to ensuring individual perpetrators are held accountable for their actions, its detractors suggest that its focus on the individual comes at the expense of building a systemic understanding of a repressive regime, and/or addressing the socio-economic impacts the apartheid state had on collectivities as well as individuals. We have identified an additional space for critique that is largely unexamined; that the volume of materials collected, combined with the TRC's focus on individual experience, obscured broader contours of violence under an archival high tide of testimonies and data collection. In a sense, the TRC translated post-Holocaust notions of 'never forgetting' as 'leaving no testimony behind' and 'no datum uncounted'. The result is an unwieldy archive that stretches across several thousand transcripts produced in archaically formatted HTML files, and a partial output of its internal database, known as *Infocomm*.

Recent developments in machine learning mark an inflection point in the field of natural language processing which has important consequences for the way we read archives and use them to write history. Although quantitative approaches to history long predate the present

moment, cliometric readings of the past largely draw from data originating in structured form (tax registers, population surveys, baptismal records, etc.). Such records were conceived of as tabular data and intended to be read as such, however advances in NLP, in particular the sub-field of named entity recognition (NER), now allows us to extract structured data from a wider variety of texts including testimony transcripts and narrative prose. Custom trained machine learning models can now automate the recognition of different categories of information (names of individuals and organizations, dates, geospatial locations, etc.) embedded in narrative texts with a high degree of accuracy. Once recognized and tagged this data can then be extracted into structured datasets that can be put to a variety of other uses. In the case of *Bitter Aloe,* it has allowed us to draw connections between 21,400 human rights violations descriptions that total over 780,000 words. In this way, *Bitter Aloe* is an attempt to overcome the unwieldiness of the TRC archive, and draw new connections within it, through the application of machine learning.

Our work began in 2019 with a partial output of data collected by TRC statement takers that was compiled into an internal database known as Infocomm. At present, only a fraction of the total data collected by statement takers and entered by TRC staff is publicly available. This fraction of Infocomm became public after South African History Archive (SAHA) prevailed in its lawsuit against the Department of Justice, leading to a judgment in 2016 that ordered the partial release of data (for the sake of clarity we will refer to the partial output of data from Infocomm as the *SAHA-Infocomm* dataset). The *SAHA-Infocomm* dataset includes only the names of individual victims, their ages, and a brief 3-4 sentence description of the human rights violations (HRVs) committed against them [include examples of descriptions]. A variety of other data, including the names of perpetrators, the gender of victims and sensitive personal information like ID numbers and street addresses were withheld. In total, some 21,399 individuals are named in this dataset, and the total length of HRV descriptions roughly equals the total amount of text in Shakespeare's corpus. Our overarching goal was to use machine learning to extract data from these descriptions and structure it into a usable form to overcome some of the limitations imposed by the partial release of the Infocomm database.

Incident descriptions entered into Infocomm originated from raw notes taken in response to questionnaires administered by statement takers as individuals came forward to tell their stories. It is unclear what became of these handwritten fieldnotes, but after collection information in the fieldnotes was condensed, abbreviated, or coded and entered into at least 60 fields within Infocomm. The varied way individuals described their experiences introduced ambiguity into the collection of data. To address this ambiguity, particularly in terms of violence used in a single act–the basic quantum of the who-did-what-to-whom approach, statement takers used a coding frame that divided human rights violations into five categories (killing, torture, severe ill treatment, abduction, associated violation) and over 63 types of violence that fell under one or more of these categories.

These details are important because they can help establish an epistemology of knowledge assumed to be represented in the publicly available *SAHA-Infocomm* dataset. Any coding frame would certainly impose dilemmas for those transcribing raw notes into structured data.

Many individuals suffered not one, but multiple forms of violence, but the brevity of descriptions often limited staff to choosing fewer than three categories. This led to an implicit prioritizing of certain forms of violence over others. The example of Selby Mavuso, an MK cadre abducted from an ANC residence during an SADF cross-border raid in January 1981. His incident description appears as follows:

An MK operative who was abducted from Matola, Mozambique, by SADF Special Forces on 30 January 1981. He was handed over to members of the Security Branch, who tried unsuccessfully to recruit him as an askari. When all attempts failed, Vlakplaas operatives took him to a spot near Komatipoort, Transvaal, where he was shot dead and his body burnt. The commander of Vlakplaas was granted amnesty for the killing, while a Vlakplaas askari was granted amnesty for his role in the attempt to recruit Mr Mavuso (AC/2000/163 and AC/2001/279).

While this is one of the longer descriptions, it foregrounds certain details. His murder by shooting and burning of his body as relevant. However, it leaves out injuries he suffered during the raid itself, the torture he endured at the hands of the Security Branch, and the fact that police attempted to induce cardiac arrest by surreptitiously poisoning his beer with an experimental chemical agent. Poisoning appears in the coding frame, but while it is unclear what information was included in the statement-taker's notes, that keyword does not appear in the description. To be sure, these descriptions are not meant to be exhaustive reports of victims' experiences. But it is important to keep in mind that a certain amount of selection, translation and foregrounding preceded the condensation and entry of incident descriptions.

TRC Vol 7 Dashboard

The main entry point into the project is a ArcGIS map called the *TRC Vol 7 Dashboard*, which maps datasets extracted from the *SAHA-Infocomm* dataset and provides a series of filters users can use to drill down to more finite research questions.¹ This work was largely the product of Robert Vaughan, our ArcGIS team leader. The first map tab titled 'cluster' represents the density of HRVs in a given location, indicated by the size of the plot for a particular locality or region. The second map tab titled 'HRV Geo' allows users to see individual HRV violations as well as the corresponding descriptions as they appear in the *SAHA-Infocomm* dataset. The third map 'Rec/Pop' normalizes the total number of HRVs in a particular municipality over 1996 population data and displays results as a heatmap of the per capita rate of violations. Each map tab can be filtered using different criteria; by pre-1994 province/homeland, HRV type, organizations referenced in incident descriptions, and year of incident. This map is a work-in-progress and we intend to add additional features and refinements.

Perhaps the most compelling feature of the map is in its ability to draw attention to regions that experienced violence at a higher per capita basis than major metropolitan

¹ Note: we initially used 'Vol 7' as a shorthand for the *SAHA-Infocomm* dataset, since Volume 7 of the TRC Final Report had been the most public-facing output from Infocomm prior to the 2016 judgment and release.

centers with higher overall HRV totals. Viewing the 'Rec/Pop' tab allows users to quickly identify regions outside major metropolitan areas that experienced anomalous rates of violence. The most disproportionate levels of violence occurred in the Ndwedwe Municipality in the Natal Midlands region, which given even a cursory read through public hearings on violence in Natal should not come as any surprise. However, more peripheral regions with smaller populations such as Ikwezi municipality in the Eastern Cape, and Langeberg municipality in the Western Cape experienced higher than average per capita rates of violence, but with less representation in public hearings, reporting and scholarship. The map offers few clues as to why violence was so acute in these rural areas, but such findings point to one use for our tools which is identifying communities and the stories within them that beg further research in conventional archives.

Anomalies revealed through mapping also point scholarly attention to the internal processes of the TRC itself, particularly in regard to the composition of public hearings. As Tutu himself states in his memoir, public hearings needed to be representative to some degree of the overall character and pattern of gross human rights abuses, but also be sensitive to the unique ways particular categories of victims experienced violence, as well as making room for what he termed 'the little people', or those individuals whose stories were unlikely to be documented and received little prior attention. There is an opacity to the selection process for public witnesses, particularly witnesses that appeared before the Human Rights Violations hearings. Understanding this opacity is important because public hearings largely shaped popular perceptions of what the TRC was and what it intended to achieve, and have by and large preoccupied scholarly debates over its significance. The disproportionate influence public hearings have had on the formation of opinions about the TRC, is belied by the fact that the hearings themselves were managed in the sense that the selection of individuals was not a random sample but an instrumental production that was conscious of public perceptions about the purpose and the legitimacy of the commission. There was an understandable bias toward foregrounding the voices of those who experienced highly publicized incidents of large-scale violence, but given the finite nature of public hearings these big events often crowded out stories from rural areas, the Bantustans and other more peripheral regions in South Africa.

The map is limited in two ways. First, we do not have address level data for incident descriptions. So the location of individual HRVs often corresponds to the geographic center of a particular locality as provided in official gazetteers. This may create a false sense of geographic specificity or conversely suggest that an incorrect location was assigned to a particular location. Along this line place descriptors could be vague, particularly in regions with contested geographies or naming conventions, changes in spelling pre/post 1994, or due to the vagaries of memory of the individual providing information to a TRC statement taker. This was particularly true in the Natal Midlands, where the scale of violence and dislocation created the kind of imprecision that defy the neat categories of databases. We did a manual validation of about 800 locations referenced that did not correspond with a lookup in available gazetteers. This proved to

be a very difficult and time-consuming task that required a lot of judgment calls about where to geographically place an HRV that was described in vague terms. Phrasing like in 'in a forest' or 'the Tugela River' created a real dilemma as we did not want to leave incidents off the map but also did not want to create an impression of false specificity. In the end we assigned a three point scale that described our degree of certainty about geographically imprecise incidents that involved judgment calls. One future task is to color code locations using this certainty scale to flag HRVs that may be placed in a general area.

Second, although the amount of data represented here is significant, it is important to keep in mind that this is not a map of all HRVs during this period. This point may seem like a truism, but there were important limits to the collection of data that prevented all who experienced violence from coming forward. First, the TRC was limited in its capacity to notify communities that statements were being taken, and individuals, particularly in rural areas, had a limited window of time in which to meet with statement takers. Second, some victims had valid skepticism about the TRC itself, particularly in regards to its amnesty process, while others stayed away because the trauma of recounting painful events paled in comparison with the potential benefit of telling one's story. Although the dataset is large enough to be considered to be statistically relevant, one must always keep in mind that on some level it is as much a map of participation as it 'as complete a picture of what happened' as mandated in the PNURA.

Co-occurrence Network Graph (CNG)

The co-occurrence network graph is a massive network graph that connects the names of individual victims and the names of organizations named in the HRV descriptions and extracted via our NER models. The graph is primarily navigable through a zoom function, but data filters also allow certain elements to be shown and hidden according to user preferences. Additionally with the data brush function portions of the graph can be manually moved for clarity. It is important to note at the outset that the CNG simply shows that an organization or organizations are named in an individual's HRV description. The CNG does not attribute responsibility to any organization for the HRV itself. Nor does the *SAHA-Infocomm* dataset include the names of any individual perpetrators. In the future we will be using joint relation recognition to extract data about the action of unnamed individuals and organizations on victims as contained in HRV descriptions. We will also be adding a chronological dimension to the network graph, but at present it is a graph of all HRVs over the term of reference for the TRC (1960-1994/5). Individuals whose descriptions do not mention an organization are displayed peripherally around the central network graph.

The network graph most clearly shows the centrality of organizations which can visually determine the relative proportionality of violence inflicted or received on individual members of those organizations. As would be expected, members of prominent political organizations such as the ANC or Inkatha have a high degree of centrality as victims of

violence committed by the opposing organization. Likewise, state security forces exhibit a high degree of centrality for the same reasons.

The real strength of the graph, however, comes in filtering data and isolating patterns in lesser understood organizations, particularly among the myriad vigilante groups that operated during the 1980s and 1990s. This method involves scanning the list for vigilante groups of potential interest, then entering their name as a string in the 'edge:target' attribute. When filtered and standalone nodes are hidden, a user can begin to parse patterns in the geographic and temporal dimensions of individual vigilante groups and begin to build a more granular understanding of their broader impact on the practice of violence.

The disambiguation of vigilante groups is also a potential strategy for exploring new research angles. A case in point is what we call the tale of the two A-Teams. Early in our validation of organization names in our dataset we came across a vigilante group known as the A-Team active in the Chesterville township outside Durban.² With the network graph we isolated nodes associated with the A-Team. However, in addition to the expected clusters of victims in Chesterville, we identified a separate smaller cluster of victims Thabong outside Welkom. Upon further investigation in the archives, we learned that these were two 'A-Teams' operating in two separate locations. This led to two questions; (1) how and why did these organizations receive their name? and (2) were the similar names evidence of collaboration with the police?

The first question begs a detour into the cultural semiotics of township violence. Newspaper reports and the TRC Special Report news digest both attribute the name of the Chesterville 'A-Team' to the popular American television show of the same name then broadcast on SABC. The original 'A-Team' was about a band of special forces Vietnam veterans wrongly accused of a crime who then go 'underground' to escape capture and right wrongs. Tracing the first attribution of the 'A-Team' to these vigilante groups may be impossible, but the selection of this symbol over others begs speculation about what the name might have meant to various parties in this conflict, and by extension open an avenue into tracing the cultural semiotics of other vigilante groups by contextualizing their cultural origin in practices of violence. To the community, a reference to the 'A-Team' may have been a critique of the racial politics of black collaboration with authorities. A central theme of the show was the relationship between a black character, B.A. Baracus whose stereotypical brute strength was tempered and redirected by the white leader of the group, Col. John "Hannibal" Smith. Was the name given by the community as an ironic parody of racial tropes and power relations in the television show? Vigilantes themselves justified their actions by portraying themselves as underdog defenders of the defenseless. In this regard, was the name self-applied and an earnest embrace of the action hero metanarrative of 80s television? Assuming, as many in the community did, that the A-Team was closely collaborating with police, another possibility is that the name was the product of a government strategy of

² TRC Special Report, Episode 20, Section 5, 29:30, https://youtu.be/TfJjTdbF-Mc

psychological warfare, or an unintended homonym with an anodyne bureaucratic nomenclature, or a derisive moniker applied by police handlers in an ironic fashion.

Returning to the second question; does it mean that there were two A-Teams in two separate theaters of struggle? Both communities in Chesterville and Thabong long suspected police complicity in the activities of the A-Teams in their area. Given the level of coordination of counterinsurgency doctrine by police in different regions, is it possible that these groups, while distinct, were of a piece? Further investigation is necessary to answer this question, but this example demonstrates the multi-directionality of lateral thinking made possible with this sort of legibility through data visualization. Although few questions can be definitively answered with machine learning and data visualization, in most instances our tools provide a new way of surveying the landscape of evidence, and in almost every case, direct users back to the archive to follow up on discovered clues.

One limitation of the graph stems from the inconsistent way data entry was conducted in Infocomm. TRC staff did not appear to have a coding frame for organizations like they did for types of violence. Consequently, a single organization can be referred to by three or more names. This problem is most pronounced with 'the police'. Although SAP is the most common term used, police are also alternately referred to as the "police", "Special Branch", "Security Branch", "security forces", or an arrest simply uses the passive voice (i.e. 'was placed under arrest...'). Although inconsistent data entry contributed to this problem, the duplication of names also stems from other sources; confusion, hearsay, and trauma impacting the memory of victims/witnesses, and the elaborate organizational structure of state security. To some degree these problems can be resolved in the graph by using inclusive filters to include all variants of names for the same organization, but in other important ways they point to the irresolution of memory, and in turn gesture to the limits of the 'who did what to whom' model of truth collection.

Text Analysis App

The beta version of our Text Analysis App went live on Wednesday, however the app began with a basic sentence embedding function focused on the *SAHA-Infocomm* dataset which we now expanded to include our new dataset drawn from TRC public hearing transcripts, which have now been cleaned and structured from large scale analysis. For the purposes of this presentation, we will limit our description to the sentence embedding function.

Sentence embedding is a method of rendering entire sentences as numerical values that are related to all other sentences in a particular corpus in a multi-dimensional vector space. The objective of sentence embedding is to use machine learning to detect deep similarities between sentences across large corpora. This ability allows users to conduct 'fuzzy' searches that are driven more by the semantic meaning, rather than matching the same keyword across multiple sentences. The advantage of this approach is it simultaneously casts a wider net with a specific semantic 'mesh' designed 'catch' sentences that convey meanings that share some basic conceptual similarity but do not

necessarily contain the same keywords. Recent applications of sentence embedding include collation of sensory experience in Holocaust testimonies and the Odeuropa Project which uses this method to explore the olfactory heritage of Europe documented across large corpora.

To use the sentence embedding function a user must first browse the *SAHA-Infocomm* dataset for an incident description of interest. Once that description is identified the user enters the ID for that particular individual into the search box to the left. Results are returned under the dropdown menu at the top of the *SAHA-Infocomm* dataset and are ranked according to their similarity.

A sentence embedding search using the incident description for Sinqokwana Ernest Malgas (ID: 7010) provides a useful example for interpreting results.³

Searching for Similarity to:

Victim: MALGAS, Sinqokwana Ernest (7010)

Description: An ANC member who was imprisoned on Robben Island from 1963 to 1977, and was subsequently detained and tortured on several occasions between 1977 and 1989 in Port Elizabeth and Cradock, Cape.

Result 1

Victim: NDABEZITHA, Joseph (13959)

An ANC member from Guguletu, Cape Town, who was detained in 1960 and again in 1963. After his 1963 detention he was sentenced to six years' imprisonment on Robben Island but was released in 1965 on appeal. He was placed under house arrest for five years. He was again detained in 1977 and tortured by members of the Security Branch. Degree of Similarity: 0.7999788522720337

Result 2

Victim: NELANI, Zoyisile William (14506)

An ANC member who was detained and tortured in 1960 and 1963 in the Transkei area, Cape. In 1971 he was again detained and tortured in John Vorster Square police station, Johannesburg by named Security Policemen. In 1976 and 1977 he was detained in the Transkei, each time for six months, during anti-independence protests. After being detained in 1979 in Umtata, Cape, he was severely tortured in an East London police station. He was charged with treason, convicted, and served five years' imprisonment on Robben Island. Degree of Similarity: 0.7597696781158447

Result 3

³ Malgas joined the ANCYL in the early 1960s, received military training in Ethiopia in 1962 and was captured enroute back to South Africa via Rhodesia in 1963, after which he spent 14 years on Robben Island. He was among the first twelve individuals to testify during the TRC's first human rights violation hearings.

Victim: QUMBELO, Mountain (17530)

A member of the ANC underground in Cape Town who was detained for many months and severely tortured by being electrocuted and suffocated by named members of the Security Police in October 1963 in Pretoria. He was given a six-year prison sentence that he won on appeal. He was then placed under a five-year banning order. In 1977 he was again detained and served a prison term on Robben Island from 1978 to 1983. Thereafter, as a UDF activist, he was repeatedly detained under emergency regulations from 1985 to 1989. See police brutality

Degree of Similarity: 0.7576255798339844

Result 4

Victim: PHUNGULA, Helia (17312)

Was detained twice - once in March 1976 and again in July 1977 - by named members of the Security Branch in Durban. In detention, he was tortured, held in solitary confinement for eight months and interrogated about his alleged involvement in political activities in the Durban area. Degree of Similarity: 0.7485864758491516

Result 5

Victim: CHOMA, Sydney Sekwati (1230) An ANC member who was detained and tortured in February 1977 in Sekhukhuneland, Lebowa, and in Groblersdal, Transvaal. In November 1979 he was found guilty of high treason in Pietermaritzburg and sentenced to 16 years' imprisonment on Robben Island. Degree of Similarity: 0.7373364567756653

Result 6

Victim: GQOLA, Frederick Bafana Makara (2886) An ANC member in the Transkei who was imprisoned for nine years on Robben Island, Cape Town, from 17 December 1964. He was again detained and tortured in 1980 and 1986. On 27 January 1987, Mr Gqola was detained on suspicion of harbouring ANC operatives and was tortured by Transkei police while in detention at the Norwood security offices in Umtata.

Degree of Similarity: 0.7371673583984375

Result 7

Victim: TSHITAHE, Ntsumbedzeni (20334) Was arrested in 1977, detained for a long period, and imprisoned for ten years on Robben Island, for his involvement in student politics. Degree of Similarity: 0.7360461354255676

Result 8 Victim: NGOBESE, Sithembiso Ernest (15190) An ANC member who was detained in Durban on 7 December 1977 and held in solitary confinement for six months in terms of the Terrorism Act. He was severely tortured, and hospitalised as a result. The Supreme Court issued a restraining order prohibiting the Security Branch from continuing to assault him. In December 1979, he was convicted of political offenses and sentenced to five years, which he served on Robben Island. In January 1987, he was again detained under the Internal Security Act, and charged but acquitted Degree of Similarity: 0.7342235445976257

Result 9

Victim: MARRAND, Wellington Thembinkosi (7585)

Was detained and tortured by members of the SAP in Durban, in March 1984. He was convicted with five others and sentenced to five years` imprisonment on Robben Island in December 1984 for recruiting people to undergo military training with the ANC.

Degree of Similarity: 0.7316596508026123

Result 10

Victim: NJAMELA, Felinyanisa Abner (15525)

An ANC member who was detained in Guguletu, Cape Town, in 1960. He was again detained in 1963 and severely tortured by named members of the police and Security Branch in Robertson, Cape. In 1967 he was detained, charged but acquitted.

Degree of Similarity: 0.7280595302581787

The five central details of Malgas' description include (A) his arrest in 1963, (B) his imprisonment on Robben Island between 1963 and 1977, (C) his release in 1977, (D) his continued detention and torture following release (E) and a geographic focus on the Cape Province. Most results match at least three details, with the first three results matching almost all five, giving a researcher interested in this case profile a decent start finding individuals with experiences similar to Malgas. However, a keyword searches for "Robben Island" returns 137 descriptions, "1977" return 170, "torture" 1,415, and so on. Even using Boolean operators identifying similar cases would be a time-consuming task with far more false positives.

That said, sentence embedding is not an exact science, leaves much room for interpretation of shared meaning and what is (and is not) an outlier in a set of results. In this regard, there is always the danger of the 'Texas sharpshooter fallacy' where a bullseye is retroactively drawn around a random cluster of shots, essentially intuiting meaning where none might be found. But in certain cases, sentence embeddings can yield 'needle in the haystack' discoveries that allow researchers to make connections that would be difficult to discover through close readings of very large corpora.

The Future

Our intention is to drive further research in existing archives by providing tools that allow researchers to ask new questions of materials collected by the TRC. Data derived from machine learning models allow users to view patterns extending across entire archives in a single frame and make semantic comparisons within large corpora that would be beyond the capacities of an individual reader. The data visualizations possible with these methods, to some extent, democratize the archive to an extent by reducing the amount of labor required to engage with historical sources through close reading alone. And while digital methods are not a substitute for close reading, with a bit of experimentation we can turn our attention to phenomena and stories not readily apparent in the large and unwieldy archives produced by truth commissions.

All that said, one point that we have debated as we have debuted tools is to what extent distant reading of human rights archives improves public discourse or creates shortcuts, opportunities for misinterpretation, or even the weaponization of findings, particularly within social media spaces. These developments in machine learning are also coinciding with the reopening of several cases from the apartheid era, each of which raise important questions about the ethical obligations that accompany making big data legible within a single and publicly accessible frame.

In the final analysis, these concerns merit further discussion, but ultimately do not obviate the public's right to know, nor place anything into the public domain that was not already available albeit in a less accessible form. Turning to concerns often voiced in Holocaust studies, advances in machine learning come at a particularly opportune time to inform public debate about human rights violations under apartheid, as many older victims have already passed on, and events begin to fade from living memory.

Lastly, as imperfect as the TRC process was, and as limited as its archive may be, there is no comparable born-digital archive that documents human rights violations at this level and scope with this level of public access. One presently intangible benefit of the application of machine learning to a large human rights archive like this may come to fruition in the conceptualization of future data-driven truth commissions. There were ways the TRC could have collected data and testimonies differently that would have had a profound impact on our understanding of past abuses. With the new forms of legibility offered by machine learning it may be possible to create a truth commission that collects individual testimonies outside of a 'who-did-what-to-whom' approach, and begins to piece together systemic understandings of violations that allow us to move more easily between individual and the collective within the same archive.