

The Thick Machine

Anthropological AI Between Explanation and Explication

Munk, Anders Kristian; Knudsen, Asger Gehrt; Jacomy, Mathieu

Published in:
Big Data & Society

DOI (link to publication from Publisher):
[10.1177/20539517211069891](https://doi.org/10.1177/20539517211069891)

Creative Commons License
CC BY-NC 4.0

Publication date:
2022

Document Version
Publisher's PDF, also known as Version of record

[Link to publication from Aalborg University](#)

Citation for published version (APA):

Munk, A. K., Knudsen, A. G., & Jacomy, M. (2022). The Thick Machine: Anthropological AI Between Explanation and Explication. *Big Data & Society*, 9(1), 1-14. <https://doi.org/10.1177/20539517211069891>

General rights

Copyright and moral rights for the publications made accessible in the public portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

- Users may download and print one copy of any publication from the public portal for the purpose of private study or research.
- You may not further distribute the material or use it for any profit-making activity or commercial gain
- You may freely distribute the URL identifying the publication in the public portal -

Take down policy

If you believe that this document breaches copyright please contact us at vbn@aub.aau.dk providing details, and we will remove access to the work immediately and investigate your claim.

The Thick Machine: Anthropological AI between explanation and explication

Anders Kristian Munk¹ , Asger Gehrt Olesen¹
and Mathieu Jacomy¹

Big Data & Society
January–June: 1–14
© The Author(s) 2022
Article reuse guidelines:
sagepub.com/journals-permissions
DOI: 10.1177/20539517211069891
journals.sagepub.com/home/bds

Abstract

According to Clifford Geertz, the purpose of anthropology is not to explain culture but to explicate it. That should cause us to rethink our relationship with machine learning. It is, we contend, perfectly possible that machine learning algorithms, which are unable to explain, and could even be unexplainable themselves, can still be of critical use in a process of explication. Thus, we report on an experiment with anthropological AI. From a dataset of 175K Facebook comments, we trained a neural network to predict the emoji reaction associated with a comment and asked a group of human players to compete against the machine. We show that a) the machine can reach the same (poor) accuracy as the players (51%), b) it fails in roughly the same ways as the players, and c) easily predictable emoji reactions tend to reflect unambiguous situations where interpretation is easy. We therefore repurpose the failures of the neural network to point us to deeper and more ambiguous situations where interpretation is hard and explication becomes both necessary and interesting. We use this experiment as a point of departure for discussing how experiences from anthropology, and in particular the tension between formalist ethnoscience and interpretive thick description, might contribute to debates about explainable AI.

Keywords

Thick description, machine learning, Clifford Geertz, computational anthropology, ethnoscience, explainable AI

This article is a part of special theme on Machine Anthropology. To see a full list of all articles in this special theme, please click here: <https://journals.sagepub.com/page/bds/collections/machineanthropology>

Clifford Geertz famously took cultural anthropology to be an interpretative science in search of meaning rather than an explanatory one in search of law (Geertz, 1973). If there is to be computational anthropology, as this special issue proposes, the Geertzian aspiration to hermeneutics might constitute a Turing test of sorts. Conceived by Alan Turing as an imitation game (Turing, 1950), the test is considered successful if a computer exhibits intelligent behavior in a way that is indistinguishable from that of a human. In the Geertzian version, it would be successful if the computer could contribute to interpreting cultural expressions in a way that is indistinguishable from the way an interpretative anthropologist would do it. As a starting point for our reflection, we ask if it is possible to play a version of the imitation game where a machine takes the place of a human when it comes to *explicating*, as Geertz puts it, “social expressions on their surface enigmatical” (Geertz, 1973: 5). Recognizing upfront that actual

computational thick description, where the machine would write interpretations of cultural texts, is not currently feasible, we settle for a version where the machine is doing the initial work of identifying deeper situations where interpretation is not straightforward and, thus, worthwhile.

We go back to Geertz, rather than newer examples of ethnographic practice, because he represents one of the most clearly voiced oppositions to algorithmic ambitions in anthropology. Ethnography has multiple disciplinary homes, but we want to revisit the last time anthropologically

¹The Techno-Anthropology Lab, Department of Culture and Learning, Aalborg University Copenhagen, Copenhagen, Denmark

Corresponding author:

Anders Kristian Munk, Learning and Culture, Aalborg University, A. C. Meyers Vænge 15, Copenhagen 2450, Denmark.
Email: anderskm@hum.aau.dk



trained ethnographers were seriously considering writing algorithms in order to model, and thus explain, culture, while being directly opposed by their equally anthropologically trained peers. By the 1960s, ethnoscientists were already experimenting with computation in their search for formal cultural laws (Hymes, 1965). That tradition, although the center of a lively controversy in its time, has today largely been purged from the ethnographic canon, not least at the hands of Geertz himself, who was one of its staunchest critics (Seaver, 2015). One of the effects of this is that mainstream ethnography now routinely defines itself as a self-evident antithesis to any kind of formalist or quantitative analysis. The 2017 Routledge Companion to Digital Ethnography (Hjorth et al., 2017) makes no mention of machine learning, artificial intelligence, or neural networks. In the opening chapter, Mike Fortun, Kim Fortun, and George E. Marcus call for a rekindling of anthropological interest in computation and make direct reference back to Del Hymes and the computational experiments of the ethnoscientists (Fortun et al., 2017). Yet the rest of the volume never fully embraces what is currently going on in the computational field. We find that surprising since, as Michael Agar puts it, “traditional social science is on the lookout for *variables*” while “ethnographers are on the lookout for *patterns*” (Agar, 2006), and so are the current generation of machine learning algorithms. Variables, in the way Agar talks about them, were for earlier generations of more rule-based AI.

You could argue, as Nick Seaver has done (2015), that the big data paradigm, i.e. the belief in the potentials of data science to revolutionize a range of academic fields, is so reminiscent of the old ambitions of the formalists that it restrains anthropology from properly rekindling its interest in computation. Even when there is ethnographic enthusiasm for complementing big data analysis by filling its gaps with thick data analysis (Wang, 2013), it is typically premised on the fact that ethnography has sought out and identified with what Seaver calls “the analogue slot” (Seaver, 2018) so as to distance itself from the reductionism of formalist analysis. After computational ethnoscience became framed as antithetical to “real” ethnography in the anthropological tradition, engagements with computation often amount to studying data science ethnographically (e.g. Boellstorff, 2013; Fisch, 2018; Grommé et al., 2018; Mackenzie, 2017; Wilf, 2013; Williams, 2018), or complementing data science with ethnography proper in a mixed methods effort to have the best of both worlds (e.g. Blok and Pedersen, 2014; Ford, 2017; Geiger and Ribes, 2011; Munk and Ellern, 2015).

In some cases, however, ethnographers have begun embracing computation in ways that do not maintain such a clear separation between big and thick. When Anne Beaulieu talks about “computational ethnography,” she explicitly describes a practice in which “code as an instrumental form of language, and computation (rather than representation or interaction) becomes much more prominent. Algorithmic interactions, calculation and generative

potential are prominent, with consequences for being in the field and for the objects of ethnographic inquiry” (Beaulieu, 2017: 34). When Dawn Nafus uses sensor data to elicit feedback from her informants (Nafus, 2018) or when Wendy Hsu uses digital methods to augment her ethnography of sound-based cultures (Hsu, 2014), computation becomes integral to fieldwork itself (see also Anderson et al., 2009). Rachel Shadoan and Alicia Dudek have argued for the use of data visualizations as what they call “scaffolding for ethnographic insights” (Shadoan and Dudek, 2013), and on the far fringes of anthropology, culturomics has emerged as an entirely computational alternative to ethnography (Leetaru, 2011; Bohannon, 2011).

It seems fair to say, then, that not all ethnographers are burdened by Geertz’ admonition against computation—largely, probably, because they do not feel constrained to a narrow definition of ethnography as “thick description and ethnographers [as] those who are doing the describing” (Geertz, 1973: 16). If the hermeneutic explication of meaning is not the only goal, there are plenty of opportunities to engage with data science and computation in ways that both expand and transform how we think of ethnography.

On the other hand, to the extent that ethnography is still, at least in part, about thick description, and in an anthropological context such as the one offered by this special issue few would dispute that it is, the old opposition to formalist ethnoscience is still tangible in the lack of computational experiments with what we have called “algorithmic sensemaking” (Munk, 2019), i.e. the use of computational techniques for hermeneutic purposes (see Nelson, 2020 for another exception). Computational thick description remains elusive, although computational anthropology (see also Bornakke and Due, 2018; Breslin et al., 2020; Munk et al., 2016) is beginning to come into its own. It would seem obvious that one of the reasons for this elusiveness is Geertz’ own denouncement of ethnoscientific algorithms in his foundational text on thick description. When he refuses to take anthropology to be an explanatory science in search of laws, he refers explicitly to formalist analysis and its attempt to reveal the cultural rules that would allow one to algorithmically “pass for a native” (Geertz, 1973: 11), such as the conventions for ordering a drink (Frake, 1964) or talking about your kin (Goodenough, 1965) in a foreign culture.

This antagonism between interpretative and explanatory ambitions is still at play in anthropological engagements with computation today, where it locates the use of algorithms solidly, but erroneously, with the latter and as antithetical to the former.

The problem with understanding the current big data paradigm as a revival of the old formalist ambitions, only with better data, is that explainability and rule-following are not at the core of the machine learning algorithms we know today. In fact, explainability is frequently debated as a central problem (Bechmann and Bowker, 2019; Cardon et al., 2018; Lipton, 2018). These methods have not brought

us any closer to writing a rule-based algorithm capable of both passing for a native and, in the process, elucidating the cultural laws that allow one to do so. Therefore, it is hard to see them as the computational realization of formalist ethnoscience. Instead, we want to argue that the Geertzian critique of the formalists provides anthropology with a different and potentially more fruitful vantage point for engaging with current advances in machine learning. What if the point is not to explain but to explicate? Could a neural network help us do that? If the algorithm is not evaluated on its ability to pass for a native and thus function as a cultural law that predicts native behavior but is instead expected to help us identify social situations with deep layers of meaning where interpretation becomes difficult, how would that change our position on explainability in AI?

Both machine learners and interpretative anthropologists are in the business of attuning to native ways of ascribing meaning to a situation; both do so without a clear theory of how that happens, in the sense that machine learning is increasingly becoming unexplainable and interpretative anthropologists have always relied on a somewhat underspecified moment of clarity (also known as the “Geertzian moment”), where the field begins to make sense; both rely on immersion to accomplish their task. Neither machine learners nor interpretative anthropologists produce a set of cultural rules that allow them or others to explain anything; neither are capable of explaining their own process in precise terms. Contrary to interpretative anthropologists, machine learners are not born with an ambition to explicate anything to an external readership (explication taken here as distinctly different from explanation, namely, as the act of constructing a reading of a situation); however, given the commonalities laid out above, one could imagine the possibility that they could be of assistance in such a process. The standards by which we judge the usefulness of machine learning techniques in anthropological analysis should change quite dramatically if the objective is not, as Geertz put it, “to codify abstract regularities but to make thick description possible, not to generalize across cases but to generalize within them” (1973).

The concept of the thick machine

In his foundational text on thick description, Geertz draws on the philosopher Gilbert Ryle and his example of the wink. A thin description, says Ryle, simply sees a muscular contraction of the eyelid. A thick description, on the other hand, must reckon with multiple possible interpretations of that muscular contraction, all of which depend on the context. Is it a flirt? Is someone mocking you or telling you a secret? Is it a sign of sarcasm or of friendship? In 2016, Facebook introduced a new set of emoji-styled reactions as an alternative to the “like.” To react with a “haha” (laughing smiley), “wow” (surprised smiley), “angry” (angry smiley), “sad” (crying smiley), or “love” (heart) to

a post can demonstrably mean very different things. A “haha” reaction can, for example, imply sarcasm, desperation, or genuine amusement, all of which radically alter the analysis of that situation. “If ethnography is thick description and ethnographers those who are doing the describing,” writes Geertz, “then the determining question for any given example of it (...) is whether it sorts winks from twitches and real winks from mimicked ones” (Geertz 1973: 16). We translate that ambition into the emoji reactions of Facebook and pursue it through an analysis of a corpus comprising 25 M posts and 128 M comments on 71 K public Facebook pages geographically located in Denmark. What would it mean to ask machine learning to help us sort ironic “haha’s” from amused ones?

To sort “winks from twitches and real winks from mimicked ones” could, in principle, be construed ethnoscientifically as the ambition to uncover the rules that separate winks from twitches in certain cultural settings and thus *explain* when a twitch constitutes a wink. That is not how Geertz intended it. Rather, for him, it should be construed in the hermeneutic sense as the ambition to sort through and elucidate the overlapping layers of meaning that make cultural situations ambiguous and hard to interpret, even for those who are involved in them. For this purpose, he distinguishes between *explanation* and *explication*. To *explicate* is to do the interpretative work of making sense of cultural texts, which is what Geertz wants. In contrast, to *explain* is to discover cultural rules, which, according to him, is what his contemporary formalists attempt to do when they pursue ethnographic algorithms capable of passing for natives.

The task of sorting winks from twitches with the aid of machine learning offers an opportunity to make this difference between explanation and explication clear through a practical demonstration and to discuss how algorithms are not necessarily so clearly aligned with the formalist position anymore. Consequently, we implement and discuss three ways of designing a Turing test for computational anthropology, two of which are modeled on the formalist ambition to explain and one on the hermeneutic pursuit of explication (see Figure 1).

In the first version, which we call *the naïve ethnoscientist*, the imitation game is simply about passing for a native. If the machine can emoji react in a way that is indistinguishable from the users on Facebook, then we consider the test to have been successful.

In the second version, *the reflexive ethnoscientist*, the goal is still to react like the users on Facebook, but the imitation game is about doing so in a manner that is indistinguishable from a group of ethnographically trained humans who are playing the game against the machine. The machine succeeds at this version of the test if it can imitate the failures and successes of the players.

Finally, in the third version, dubbed *the interpretative ethnographer*, the goal is not to pass for a native but to identify situations where overlapping layers of meaning make thick

description possible and worthwhile. Doing so is arguably the first task for any hermeneutic, Geertzian-inspired, interpretative cultural analysis. This means that we should be able to sort the emoji reactions where the meaning is fairly unambiguous and obvious to the participants from the ones where several possible interpretations are going on as part of the situation. The game is no longer about imitating the users on Facebook per se but about identifying situations where interpretation becomes complicated, thus imitating the interpretative ethnographer in the early stages of fieldwork.

In the following, we recount how we built a machine that allowed us to play these three versions of the imitation game. We call it the Thick Machine, which is deliberately ambiguous in its own right. The machine turns out to be both as thick (i.e. clueless) as the players when it comes to predicting emoji reactions on Facebook, and somehow still capable of pointing us in the direction of interesting thick descriptions (i.e. explications) of those reactions. We will show how the ethnoscientific versions of the game run into well-known problems with explainability in machine learning. Ultimately, it is not enough, from the perspective of a formalist cultural analysis, to have an algorithm capable of passing for a native. It must also result in discernible cultural rules. Interestingly, the interpretative version of the imitation game turns out not to have the same problem. Displacing the ambition from explanation to explication also changes the anthropological expectations of the algorithm.

The dataset

The Thick Machine makes use of a dataset that we harvested in January 2018 to map the debate on all public Danish Facebook pages (see Munk and Olesen, 2020 for background and details). The full dataset for the Atlas of Danish Facebook Culture, as we called the project, comprises 25 M posts and 128 M comments from 71 K Danish pages in the period between January 2012 and January 2018. These posts and comments are associated with 700 M like reactions and 23 M emoji reactions. For each post and comment, we store the full text and the reactions that have been left by users in response to that text. We also store a unique user ID for the authors of each post, comment, and reaction. We use an Elastic Search database, allowing us to easily associate a comment on a post with the comment author's reaction to that post. The combination of a comment and an emoji reaction made by the same user to the same post constitutes the basic unit of analysis for the Thick Machine and the players competing against it. Knowing the text of the comment, the task for the player is to predict the correct emoji reaction.

As a consequence of the protocol we used (see Figure 3), the Atlas of Danish Facebook Culture does not cover all public Danish pages. We began by covering the territory of Denmark with 1086 geolocations. We then used the Facebook Place Search API to find 69 K pages with an

Which game?	Imitating who?	By doing what?	Criteria for winning
<i>The naive ethnoscientist</i>	Users on Facebook	Predicting the users' emoji reactions	High prediction accuracy
<i>The reflexive ethnoscientist</i>	Ethnographers trying to imitate users of Facebook	Predicting emoji reactions in the same way as the ethnographers	Similar pattern of accuracy as the ethnographers
<i>The interpretative ethnographer</i>	Ethnographers trying to identify situations on Facebook that need explication	Showing where emojis are hard to predict	Finds situations where multiple layers of meaning complicates interpretation

Figure 1. Three versions of the imitation game for a Turing test in computational anthropology.

address within a 5 km radius of one of these geolocations. Most of these pages had been automatically generated by Facebook to correspond to a landmark. They did not have an administrator, and they did not host a discussion on their wall. Only 4600 of them were genuine pages with an administrator and an active wall. We used them as the seed for a snowballing strategy, where we asked the Facebook Graph API to return pages liked by the seed pages, check if the returned pages had an address in Denmark, and, if so, add that page to the corpus. We had to adopt this snowballing strategy because the Place Search API did not return anything near a comprehensive result on the geographically delimited place searches. Instead of the 4 K pages originally returned by the Place Search API, the snowballed corpus from the Graph API ended up covering 71 K pages with a physical address in Denmark after 15 iterations of the method. We can be sure that this is not all pages in Denmark because many pages do not have a physical address at all and will therefore, by definition, be overlooked by our protocol. An alternative would have been to use a language criterion, but since we know that more than 20% of the pages in our corpus are non-Danish speaking (Munk and Olesen, 2020), this would have introduced another bias and produced another form of incompleteness. We therefore opted for the geographical criterion, and from the 71 K pages that satisfied this criterion, we harvested all posts, comments, and user reactions.

Emoji reactions (“wow,” “sad,” “angry,” “haha,” and “love”) were introduced by Facebook as an alternative to the “like” in early 2016. To construct the dataset for the Thick Machine, we therefore focused on posts and comments from 2016, 2017, and 2018. We identified comments where the author had also left an emoji reaction on the post by checking if the user ID of the commenting author was also among the user IDs reacting with an emoji to the post. This produced a list of post-comment pairs where we knew the emoji reaction to the post by the commenting author (see Figure 2). If the same author had commented several times on the same post, we selected the first comment only. In order to be able to train classifiers on the comment text and to give the players a chance of interpreting the comment in the context of a post, we set a

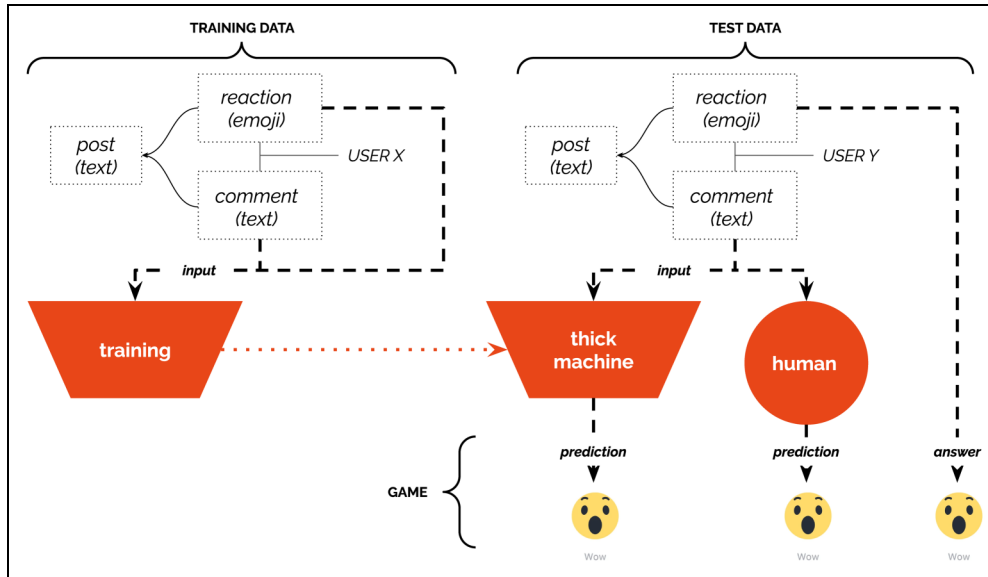


Figure 2. Blueprint for the thick machine. On the left, a neural network is trained to associate emoji reactions with comment texts. When the same user comments AND reacts with an emoji on the same post, the comment text and the emoji reaction are associated and fed to the machine for training. On the right, the trained machine tries to predict emoji reactions for comments it has never seen before, as does a group of human players who compete against the machine. Their predictions are then evaluated against the actual emoji reaction used by the author of the comment in question.

minimum text length of 50 characters for both posts and comments. As we shall see in the next section, the classifiers we use require the comment text to be broken down into individual words and those words to be transformed into numbers that describe their relative significance, i.e. a feature space. Longer comments provide more features for classifiers to train on. This left us with a dataset of 700 K post-comment pairs.

We knew from earlier research (Munk and Olesen, 2020) that the “love” reaction was the most prevalent in public Danish Facebook discourse. This produces a bias in the dataset (Figure 4, left). To ensure a 20% chance that any one of the five emoji reactions would be associated with a comment in our game, we therefore randomly deleted comments with overrepresented emoji reactions until they had the same frequency as the least prevalent (the “sad” reaction), reaching an evenly balanced dataset of 175 K post-comment pairs with 35 K pairs for each of the five reactions (Figure 4, right).

Training the machine

In order to train the machine, we first performed feature extraction, comparing a term frequency-inverse document frequency (TF-IDF) to a latent Dirichlet allocation (LDA) approach. We then split the data into a training set and a test set in order to let a neural network learn how to classify emoji reactions associated with comments in the training set and subsequently predict emoji reactions for comments in the test set. Here, we also experimented with different

approaches, eventually settling on a combination of TF-IDF-based feature extraction and a multilayer perceptron (MLP) classifier with two hidden layers trained on 80% of the dataset. This produced 51% accuracy when predicting emoji reactions in the test set.

Feature extraction is the process of assigning weight to each word in a comment text. The classifier cannot be fed raw text but requires a feature space. In this input table, the rows are the comments from Facebook, and the columns are the features, in this case, words found in the comments. The task of feature extraction is to decide how important each word is in each of the comments from Facebook and assign a numeric weight accordingly. The most basic approach would be to simply decide whether a word is in a comment or not. Some words, however, are better at distinguishing the particular tone or content of a comment. We wanted the classifier to place more weight on these words than on others. Therefore, it became necessary to evaluate the importance of a word in more sophisticated ways.

Deciding how to weight words according to their significance in a text is a fairly routine task in natural language processing (NLP), and there are multiple possible approaches. We tried two different “bag of words” models, so-called because they disregard the order of words in the text. Each word in the set becomes a feature, eventually producing a vast number of features (the dimensions of the feature space) when iterating over many documents. The question is how to rank them. Our first strategy was aimed at simplicity. Using the TF-IDF metric, we counted how many times a word appears in a document,

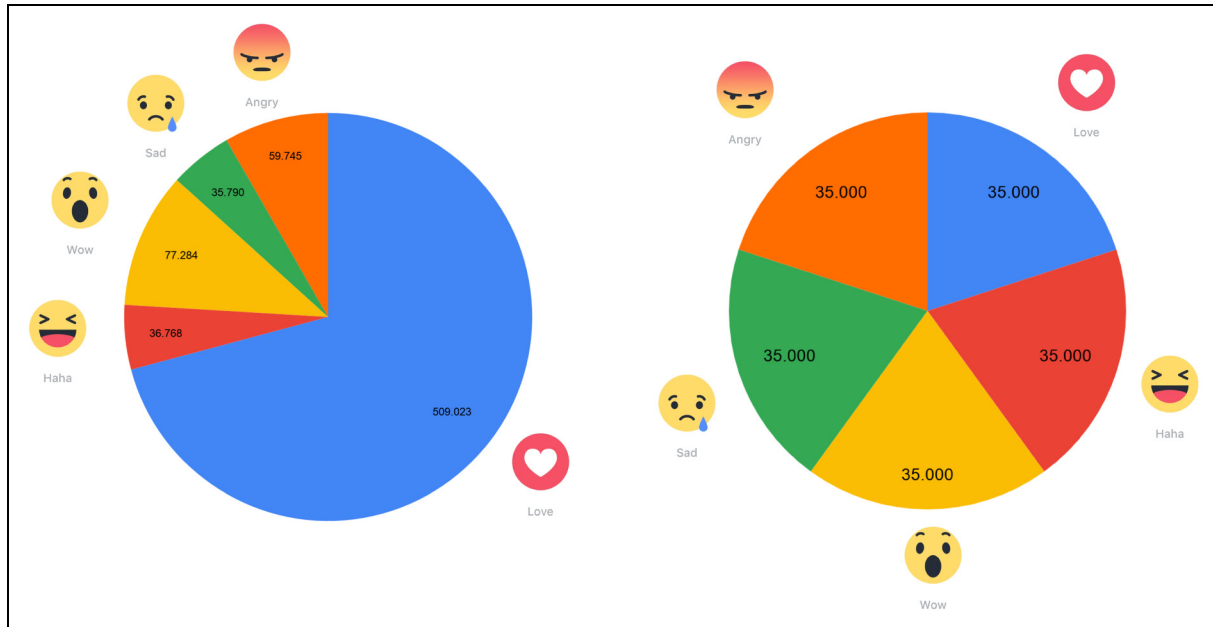


Figure 3. The protocol used to harvest and delimit the dataset for the Atlas of Danish Facebook Culture (Munk and Olesen, 2020).

normalized by its frequency across all the documents in the corpus. Our second strategy looked at groups of words. We used LDA to build groups that characterize documents (also known as topic modeling). Each document has a weighted link to each group, corresponding to how many times the words of the group are present in the document. Words that frequently appear together end up in the same group. In this strategy, the groups (or topics) become the features, taking the value of the weight of the association between the group and the document.

Once we had built the feature space, we used it to train a machine learning algorithm to classify comment text as being associated with either a “haha,” “wow,” “sad,” “love,” or “angry” reaction. Since we had a tagged training set (the feature space was already associated with the correct emoji reactions), we could use a supervised approach, where classifiers are trained to predict a known target. The goal is to be able to attribute one of several possible classes (the five emoji reactions) to a given item (a comment from Facebook). We used the library SciKit Learn to build our classifier in Python and settled on the MLPClassifier module, using interconnected layers of simple decision models inspired by the (biological) neuron. Once the number of neurons and layers is configured, the classifier learns from the training dataset by iteratively adjusting the weights of the connections between the neurons. After this process, the *trained* classifier consists of a *weighted* network of neuron connections.

As mentioned, the best accuracy we were able to achieve was 51%, two-and-a-half times better than random (20%). This result was obtained with the TF-IDF approach,

which is both simpler and less time-consuming than the LDA approach. Increasing the volume of training data significantly improved the accuracy (Figure 5). Since we did not need a large test dataset for the game itself, we set aside 80% of the data for training. Somewhat surprisingly, the number of layers in the MLP classifier did not significantly affect the accuracy. Two layers with 1000 and 100 neural nodes instead of a single layer with 100 neural nodes only increased accuracy by 2%.

Building the machine

Once we had trained the classifier, the idea was to let the researchers in our laboratory play against the machine. They are all experienced Facebook users and trained ethnographers. Since the goal was not to have a representative sample of humans from any particular population play the game but to make a practical demonstration of the difference between an ethnoscientific and a Geertzian approach to machine learning, we determined that the group of players was adequate for our purposes. The objective was to test how they would fare having to interpret the Facebook comments and predict the same emoji reactions as the classifier.

Instead of carrying out the test with a simple questionnaire on a computer screen, we decided to manifest it as a physical arcade game for two reasons. First, it seemed fair to level the playing field and put the players in a situation that was as similar as possible to that of the machine. The classifier quite literally sees text with no context. It does not have the benefit of scrolling a thread of comments, browsing through the post activity on a page, or scanning

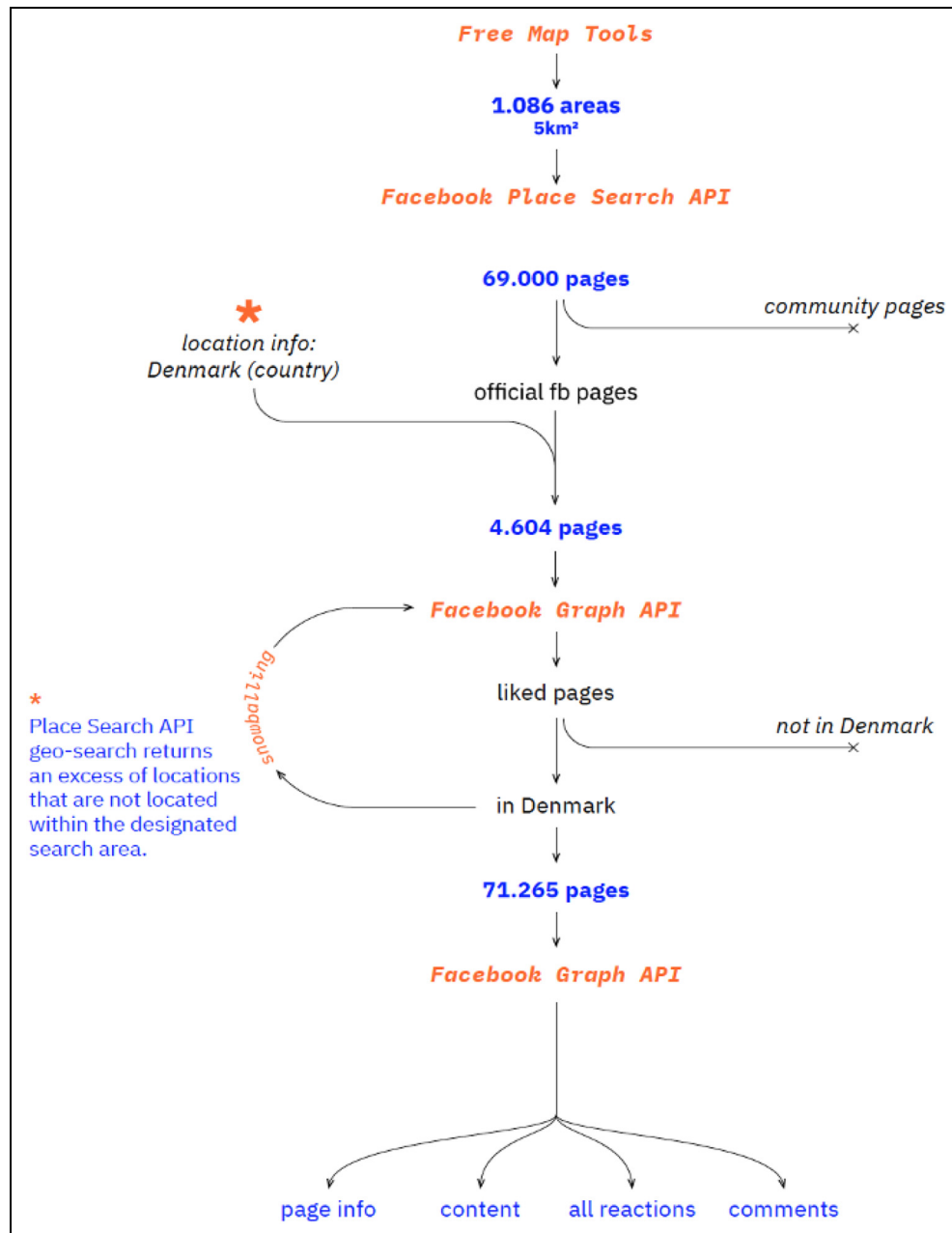


Figure 4. Randomly reducing the dataset (left) to ensure that all reaction types are equally present (right).

a user's past activity to put an interaction into context. It does not see the name or profile picture of the author of a comment and thus does not infer information about that person. This could, in principle, be implemented if the data were available and priority was given to a more complicated machine learning workflow. Indeed, one could argue that the ability to understand a Facebook comment in the wider context of what is happening on the platform should give the human player a distinct advantage. However, since the data were not available for us to give

the machine a fair chance at doing something similar, we decided instead to approximate the gaming experience to what the machine was actually training on. By designing an arcade game with its own set of clunky and unfamiliar controls (see Figure 6), and where the player is presented with raw text in a primitive green-on-black display style, the idea was to momentarily deprive the players of some of the contexts they would otherwise exploit.

Second, we wanted the game to be easily playable by the researchers in our lab when it was convenient for them.

<div>classifier</div> <div>features</div>	MLPC: 1 hidden layer, 100 nodes (Multi Layer Perceptron Classifier)	MLPC: 2 hidden layers, 1000 + 100 nodes (Multi Layer Perceptron Classifier)
TF-IDF (Term Frequency - Inverse Document Frequency)	Trained on 1% of dataset: 41% accurate Trained on 80% of dataset: 49% accurate	Trained on 1% of dataset: 42% accurate Trained on 80% of dataset: 51% accurate
LDA (Latent Dirichlet Allocation)	Trained on 1% of dataset: <i>not done</i> Trained on 80% of dataset: 31% accurate	Trained on 1% of dataset: <i>not done</i> Trained on 80% of dataset: 32% accurate

Figure 5. The results of our experiments with different approaches to feature extraction (term frequency-inverse document frequency [TF-IDF] and latent Dirichlet allocation [LDA]), different settings for the multilayer perceptron (MLP) classifier, and different proportions in the amounts of test and training data.

Building the physical arcade machine allowed us to put it on the main conference table for a period of five days. During that time, lab members made 188 predictions while playing against the Thick Machine.

Evaluating the imitation game

If the game is that of “the naive ethnoscientist” who is simply trying to pass for a native, the goal is to have as high accuracy as possible (see Figure 1). Both the players and the machine fail in about half of their predictions (see Figure 7). The machine, however, is more consistently accurate across the different emoji reactions. The players manage to predict 82% of the “love” reactions but only 25% of the “wow” reactions correctly. The machine manages to predict 56% of both the “sad” and the “angry” reactions but only 44% of the “wow” reactions correctly. The same difference is visible if we evaluate the columns rather than the rows in the confusion matrix (Figure 7). When the players predict that the

emoji reaction associated with a comment is “haha,” they are correct in 83% of cases, but when they predict that the reaction is “wow,” they are only correct in 33% of cases. Conversely, when the machine predicts that the emoji reaction associated with a comment is “angry,” it is correct in 55% of cases, but when it predicts that the reaction is “wow,” it is only right in 44% of cases. So, we can conclude that the players do better than the machine when the reaction is “love,” “angry,” or “haha,” and the machine does better than the players when the reaction is “wow” or “sad,” but both are relatively bad at imitating the users on Facebook.

If the imitation game is, instead, that of “the reflexive ethnoscientist,” where the machine is compared to players trying to imitate the users on Facebook, then we should ask the question differently. The important thing is no longer whether the machine is wrong or right but whether it is wrong and right in the same way as the players. The overall accuracy of 51% for the machine and 52% for the players means that, overall, the machine is successfully



Figure 6. Construction of the arcade machine with a Raspberry Pi, 10 emoji-styled push buttons (for two-player mode), and the PyGame interface running in the background before launching the black-on-green display of the game itself.

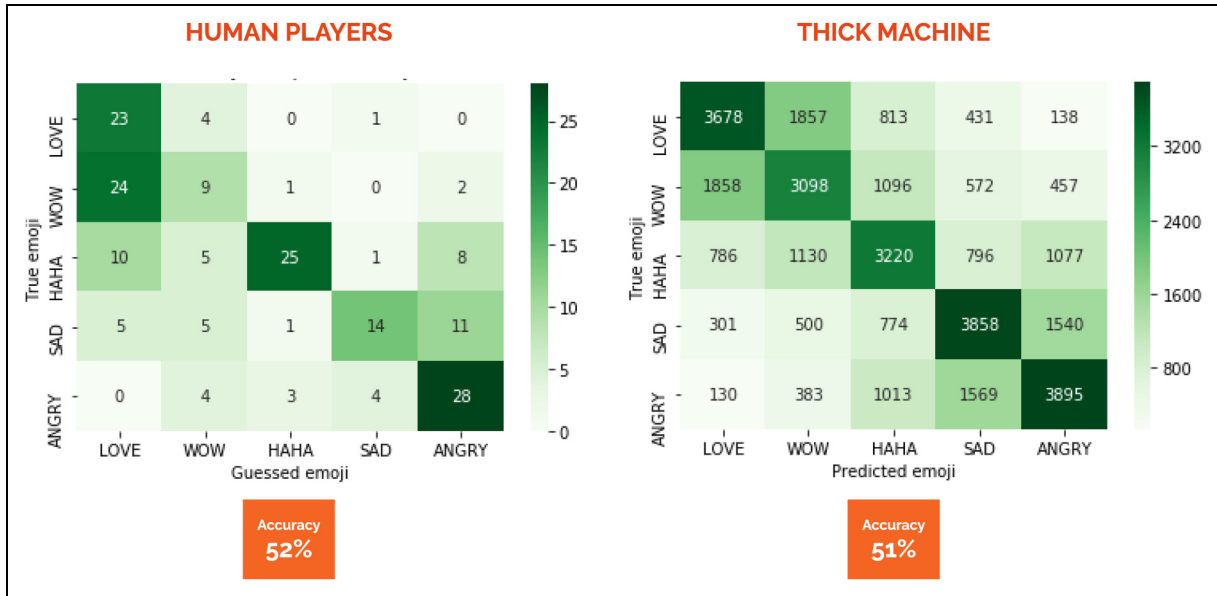


Figure 7. Confusion matrices comparing the performance of the human players in the TANTLab (left) to the performance of the neural network in the thick machine (right).

imitating the players. Again, there are notable differences in the level of each reaction type. In some instances, the players and the machine make mistakes that are remarkably similar. The reaction “angry” is predicted instead of “sad” in 22% of cases for both the players and the machine, “angry” instead of “haha” in 16% of cases for the players and 15% for the machine, “haha” instead of “love” 16% of the time for players and 12% for the machine, and “wow” instead of “haha” 19% of the time for players and 16% for the machine. To some extent, then, we can say that the machine could be successfully mistaken for a player.

In other instances, however, the machine makes mistakes that the players do not. This is particularly true for “wow,” which is mistaken for “haha” 16% of the time by the machine against 3% by the players, and “love,” which is mistaken for “haha” 12% of the time by the machine against 0% by the players. And then there are the instances where neither the players nor the machine tends to make mistakes. This is especially the case for the “love” and “angry” reactions, which are almost never confused. The players never predict “angry” when the correct reaction is “love,” nor do they predict “love” when the correct reaction is “angry.” The machine makes those mistakes in only 2% of cases. Since “love” and “angry” are arguably the two most mutually exclusive of any pair of emoji reactions (contrary to “sad” and “angry” or “wow” and “angry,” e.g.), it is perhaps not surprising that neither players nor machine tend to confuse them.

Finally, the players generally overestimate how much love there is in the dataset. When they predict “love,” they are only correct in 37% of cases, pulling their overall accuracy down. The machine does not make that mistake as often and predicts “love” correctly in 54% of cases. This could

likely be due to our normalization of the dataset to balance the frequency of the five emoji reactions. If the players’ experience approximates the average Danish Facebook user, they will be expecting to encounter the “love” reaction somewhere between 7 and 14 times more frequently, depending on which of the other emoji reactions we are comparing it to. The machine, in comparison, only knows the balanced dataset and does not suffer from this bias.

If we were playing the “reflexive ethnoscientist” version of the imitation game, where the machine tries to predict and make mistakes in a way that resembles the players, we could be satisfied to have had a partly successful test here. The imitation works when it comes to confusing “sad” reactions for “angry” reactions or never confusing “angry” reactions for “love” reactions. It still needs some improvement when it comes to distinguishing between “haha” reactions and “wow” reactions, or when it comes to estimating the overall volume of love in the set. However, it is conceivable that we would be able to gain some headway on these points by training the machine on a dataset that was not artificially balanced between the reaction types. If the objective of the game is to imitate how players predict emoji reactions, the classifier should rather be trained on a dataset where the distribution between reaction types resembles what the players will likely have experienced as ordinary Facebook users.

From explanation to explication

The key ambition of this paper, however, was to explore what happens if the concept of the imitation game is not modeled on the ethnoscientific idea of an ethnographic

algorithm capable of passing for a native. What if the game did not center on the notion that there are cultural rules that can be uncovered but on a version of Geertzian thick description? If the imitation game we are playing is that of “the interpretative ethnographer,” then it cannot be about accuracy in predictions. As we have already discussed in the introduction, Geertz objected to the idea that the measure of success for a thick description is the degree to which it mimics a native interpretation of events. In his critique of formalist ethnoscience, he thus directly addressed what he saw as a problematically algorithmic approach to the interpretation of culture:

Various called ethnoscience, componential analysis, or cognitive anthropology (a terminological wavering which reflects a deeper uncertainty), this school of thought holds that culture is composed of psychological structures by means of which individuals or groups of individuals guide their behavior. “A society’s culture,” to quote Goodenough again, this time in a passage which has become the locus classicus of the whole movement, “consists of whatever it is one has to know or believe in order to operate in a manner acceptable to its members.” And from this view of what culture is follows a view, equally assured, of what describing it is—the writing out of systematic rules, an ethnographic algorithm, which, if followed, would make it possible so to operate, to pass (physical appearance aside) for a native. (Geertz, 1973: 11)

If we take the game to be directly about imitating the natives on Facebook, or indirectly about imitating the players who are in turn trying to imitate the natives on Facebook, then we are by extension also subscribing to the formalist ambition of an ethnographic algorithm capable of passing for a native. Indeed, in those versions of the game, the failure of the machine to properly imitate either players or natives would only be the beginning of our problems. A much greater challenge would arise when we had to explain how the algorithm works. Based, as it is, on a neural network, the Thick Machine could never satisfy an ethnoscientific expectation of algorithms as something akin to the “writing out of systematic rules.”

In the debate on explainable AI (overview in Wieringa, 2020), some researchers argue that there might be a trade-off in machine learning techniques between accuracy and explainability. Explainability refers here to our ability to scrutinize how an algorithm makes its decisions and assess its biases. Decision trees, for example, are explainable by allowing us to retrace the decision-making process, but unfortunately, they are also quite inaccurate. Conversely, deep learning (i.e. multilayered neural networks like the ones we use) performs much better but is poorly explainable because we cannot easily retrace how classification is made.

Some researchers (e.g. Lipton, 2018; Liu et al., 2018) argue that algorithm accountability should not be based on explainability but on post hoc interpretability, although this argument has its detractors (Laugel et al., 2019). Post hoc interpretability assesses algorithms as black boxes for the results they produce, independently of how they function. Indeed, it is argued that it might be more productive to investigate highly complex processes as empirical phenomena from the outside:

While post hoc interpretations often do not elucidate precisely how a model works, they may nonetheless confer useful information for practitioners and end users of machine learning. Some common approaches to post hoc interpretations include natural language explanations, visualizations of learned representations or models, and explanations by example (e.g. this tumor is classified as malignant because to the model it looks a lot like these other tumors). (Lipton, 2018)

To us, there is a strong resonance with thick description here. Interpretative anthropologists typically cannot formally explain the process by which immersion and participant observation led them to be able to classify, say, some cock fights, to use Geertz’ own example, as being ripe with meaning and therefore interesting to explicate. But they can say that, after going through the ordeals of obtaining rapport with the Balinese villagers, these cock fights have acquired a special significance for them, and they can validate their classifications with their informants (Geertz, 1973: 412–454). Post hoc interpretability, in other words.

It is important to stress that we are not directly comparing the deep play of the cock fight itself to our experiments with the Thick Machine. What we are comparing is the way ethnographers account for their interpretations of situations, such as cock fights, and the way in which we account for the predictions performed by a neural network. When Geertz insists that thick description does not rely on the codification of abstract regularities and generalization across cases but on the ability to generalize within a case (Geertz, 1973), his argument is precisely that a convincing reading of the different layers of meaning in a situation does not enable us to infer a set of rules that specify how such a reading can be applied to other situations.

Cardon et al. (2018) provide an insightful perspective on the disruption caused by deep learning to formalist ambitions in machine learning. They label deep learning classifiers as inductive machines, in opposition to the more classical hypothetical-deductive machines. The technique of neural networks is not new, but it was marginalized throughout much of the history of artificial intelligence, where the dominant paradigm was symbolic (it considered thinking to be analogous to manipulating symbols). Neural networks are, on the contrary, connexionist: they “think” throughout a massive set of parallel elementary

What Geertz proposes as an alternative to explanation and rule-based algorithms that pass for natives is explication. To him, culture is already explications upon explications upon explications. The ethnographer is simply adding another layer: "Right down at the factual base, the hard rock, insofar as there is any, of the whole enterprise, we are already explicating: and worse, explicating explications" (Geertz, 1973: 9). In Geertz' eyes, the ethnoscientific aspiration for cultural algorithms wrongly assumes that there is a single unambiguous native rule or logic that prescribes how a situation should be understood and, thus, holds explanatory power. Instead, there is a play going on between different "frames of interpretation" that are "ingredient in the situation" (Geertz, 1973: 9), and it is this play that the ethnographer has to interpret or, as he puts it, construct a reading of:

If that is the challenge the interpretative anthropologist faces, then what kind of games are we really playing on the Thick Machine? Some situations are more ambiguous than others; the superimposition of conceptual structures is more complex in some situations than others, and these are the situations that interpretive anthropology is especially interested in explicating. It follows that Facebook comments where it is easy to predict the associated emoji reaction are not particularly interesting for a thick description (the meaning is already clear). Below, we have randomly selected five examples (anonymized and translated

Post: “Smart: Text a live chicken to a poor family! Text CHICKEN to 1911 (20 kr), and donate a chicken. This way you help them help themselves. Do write a comment in the

a failure to predict is more interesting than accuracy, and the lack of explainability in neural networks is less of a problem.

Conclusion

We built a machine that tries to associate Facebook comments with emoji reactions. We let researchers in our lab compete against the machine in a game that allowed us to reflect on the difference between formalist and interpretative approaches to machine learning. In a simple formalist version, the game can be about imitating users on Facebook. Neither the machine nor the players are very good at this version. In a more reflexive formalist version of the game, where the goal is to imitate the players, the machine performs better. It makes many of the same mistakes. Regardless of their differing degrees of success, both versions are modeled on the ethnoscientific ambition to discover cultural algorithms capable of passing for (by reacting like) natives in specific situations. Interestingly, they also illustrate something fundamentally incommensurable about the state of machine learning today and the kind of algorithms that the ethnoscientists had in mind when they were trying to establish a computational anthropology in the 1960s.

The formalists were fond of cultural algorithms as a way to make rules explicit. To explain culture was to discover the rules that one would have to follow in order to pass for a native. The neural networks that we have used to build the Thick Machine, and which are now pervasive in machine learning, are unexplainable by these standards and could never be used for the ethnoscientific purpose of crystalizing laws of culture, even if they reached 100% accuracy in their predictions. Understanding today's emergent computational anthropology as a revival of formalist ethnoscience thus takes us into the same maze of arguments about unexplainable AI that quantitativist computational social science is increasingly finding itself mired in.

We began this paper by asking whether it was possible to imagine a form of computational thick description, deliberately reviving the fiercest contemporary critic of the ethnoscientists and their algorithms, Clifford Geertz, and suggesting that interpretative anthropology might offer a valuable alternative to the entrenched positions in the current debate about explainability and AI. If the ambition is not explanation but explication, we have argued, it reconfigures the way computational anthropology could imagine itself and its engagements with machine learning. Although we have not actually asked a computer to do thick description for us, when we shift the game to be about the machine imitating an ethnographer in search of deep situations to explicate, three interesting things happen.

First, the situations where the machine fails come into view as situations with a real potential for thick description. When it is hard to train a classifier, it is typically because the situation is ambiguous. As soon as there is a play between multiple and overlapping frames of interpretation ingredient

in the situation, both the machine and the players fail. Arguably, these failures could be understood as indications that here are situations worthy of being explicated.

Second, it is no longer a matter of replacing the ethnographer with an algorithm (which could be the eventual result in the formalist versions of the game). The ambition is now more modest, and it becomes clearer where data science ends and ethnography begins under the umbrella of computational anthropology. The situations in which the machine is able to point out as interesting for thick description are not actually thickly described by the machine. The machine, through its failures, identifies potentially deep situations, but the interpretative work of constructing readings of those situations remains in the hands of the ethnographer.

Third, the explainability of the neural network is no longer an issue in the same way. The so-called Geertzian moments, where the field begins to make sense, were never explainable in the way we typically demand it of explainable machine learning. Anthropology has developed a methodological repertoire for thinking about and coming to terms with that fact, a repertoire that offers computational anthropology an interesting position from which to make a novel contribution to the debate on explainability in AI. Our experiments with the Thick Machine have at least demonstrated that it might be possible to use neural works in a way where post hoc interpretability is perfectly in line with existing ethnographic practices.



Declaration of conflicting interests

The author(s) declared no potential conflicts of interest with respect to the research, authorship, and/or publication of this article.

Funding

The author(s) received no financial support for the research, authorship and/or publication of this article.

ORCID iDs

Anders Kristian Munk  <https://orcid.org/0000-0002-5542-3065>
Mathieu Jacomy  <https://orcid.org/0000-0002-6417-6895>

References

- Agar M (2006) An ethnography by any other name. *Forum Qualitative Sozialforschung/Forum: Qualitative Social Research* 7: 4.
- Ahmad T, Akhtar H, Chopra A, et al. (2014, September) Satire detection from web documents using machine learning methods. In: 2014 international conference on soft computing and machine intelligence, pp. 102–105. IEEE.
- Anderson K, Nafus D, Rattenbury T, et al. (2009) Numbers have qualities too: Experiences with ethno-mining. In: *Ethnographic praxis in industry conference proceedings* (Vol. 2009, No. 1, pp. 123–140). Oxford, UK: Blackwell Publishing Ltd.
- Beaulieu A (2017) Vectors for fieldwork. In: Hjorth L, et al. (eds), *The Routledge Companion to Digital Ethnography*, 29, pp. 55–65. London: Routledge.

- Bechmann A and Bowker GC (2019) Unsupervised by any other name: Hidden layers of knowledge production in artificial intelligence on social media. *Big Data and Society* 6(1): 1–11.
- Blok A and Pedersen MA (2014) Complementary social science? Quali-quantitative Experiments in a Big Data world. *Big Data and Society* 1(2): 1–16.
- Boellstorff T (2013) Making big data, in theory. *First Monday* 18: 10.
- Bohannon J (2011) Google books, wikipedia, and the future of culturomics. *Science (New York, NY)* 331(6014): 135.
- Bornakke T and Due BL (2018) Big-thick blending: A method for mixing analytical insights from big and thick data sources. *Big Data and Society* 5(1): 1–16.
- Breslin SD, Enggaard TR, Blok A, et al. (2020) How we tweet about coronavirus, and why: A computational anthropological mapping of political attention on Danish twitter during the COVID-19 pandemic. In: The COVID-19 Forum III. University of St Andrews (Vol. 18).
- Burfoot C and Baldwin T (2009, August) Automatic satire detection: Are you having a laugh? In: Proceedings of the ACL-IJCNLP 2009 conference short papers, pp. 161–164.
- Cardon D, Cointet JP and Mazières A (2018) Neurons spike back. *Réseaux* 5: 173–220.
- Fisch M (2018) *An Anthropology of the Machine: Tokyo's Commuter Train Network*. Chicago, IL: University of Chicago Press.
- Ford H (2017) The search for wikipedia's edges. In: Hjorth L, et al. (eds) *The Routledge Companion to Digital Ethnography*, pp. 416–425. London: Routledge.
- Fortun M, Fortun K, Marcus GE (2017) Computers in/and anthropology. In: Hjorth L, et al. (eds) *The Routledge Companion to Digital Ethnography*, pp. 11–20. London: Routledge.
- Frake CO (1964) How to ask for a drink in subanun. *American Anthropologist* 66(6): 127–132.
- Geiger RS and Ribes D (2011) Trace ethnography: Following coordination through documentary practices. In: System Sciences (HICSS): 1–10. IEEE.
- Geertz C (1973) *The interpretation of cultures*. New York: Basic books.
- Goodenough WH (1965) Yankee kinship terminology: A problem in componential analysis 1. *American Anthropologist* 67(5): 259–287.
- Grommé F, Ruppert E and Cakici B (2018) Data scientists: A new faction of the transnational field of statistics. In: *Ethnography for a Data-Saturated World*, pp. 33–61. Manchester: Manchester University Press.
- Hjorth L, Horst H, Galloway A and Bell G (eds) (2017) *The Routledge Companion to Digital Ethnography*. London, Routledge: Taylor and Francis.
- Hsu WF (2014) Digital ethnography toward augmented empiricism: A new methodological framework. *Journal of Digital Humanities* 3: 1.
- Hymes DH (Ed.) (1965) *The Use of Computers in Anthropology*. (Vol. 2). Mouton: The Hague.
- Laugel T, Lesot MJ, Marsala C, et al. (2019) The dangers of post-hoc interpretability: Unjustified counterfactual explanations. arXiv preprint arXiv:1907.09294.
- Leetaru K (2011) Culturomics 2.0: Forecasting large-scale human behavior using global news media tone in time and space. *First Monday* 16: 9.
- Lipton ZC (2018) The mythos of model interpretability. *Queue* 16(3): 31–57.
- Liu N, Huang X, Li J, et al. (2018, July) On interpretation of network embedding via taxonomy induction. In: Proceedings of the 24th ACM SIGKDD international conference on knowledge discovery and data mining, pp. 1812–1820.
- Mackenzie A (2017) *Machine Learners: Archaeology of a Data Practice*. MIT Press.
- Munk AK (2019) Four styles of quali-quantitative analysis: Making sense of the new nordic food movement on the web. *Nordicom Review* 40(s1): 159–176.
- Munk AK and Ellern AB (2015) Mapping the new Nordic issue-scape: How to navigate a diffuse controversy with digital methods. In: *Tourism encounters and controversies: ontological politics of tourism development*, pp. 73–96. London: Routledge.
- Munk AK, Abildgaard MS, Birkbak A, et al. (2016, July) (Re-) Appropriating instagram for social research: Three methods for studying obesogenic environments. In: Proceedings of the 7th 2016 International Conference on Social Media and Society, pp. 1–10.
- Munk AK and Olesen AG (2020) Beyond issue publics? *Curating a Corpus of Generic Danish Debate in the Dying Days of the Facebook API*. STS Encounters-DASTS working paper series 11: 1.
- Nafus D (2018) Working ethnographically with sensor data. In: *Ethnography for a Data-Saturated World*, pp. 233–251. Manchester: Manchester University Press.
- Nelson LK (2020) Computational grounded theory: A methodological framework. *Sociological Methods and Research* 49(1): 3–42.
- Rubin VL, Conroy N, Chen Y, et al. (2016, June) Fake news or truth? Using satirical cues to detect potentially misleading news. In: Proceedings of the second workshop on computational approaches to deception detection, pp. 7–17.
- Seaver N (2015) Bastard algebra. In: Boellstorff T and Maurer B (eds) *Data, Now Bigger and Better*, 27–45. Chicago, IL: Prickly Paradigm Press.
- Seaver N (2018) What should an anthropology of algorithms do? *Cultural Anthropology* 33(3): 375–385.
- Shadoan R and Dudek A (2013) Plant wars player patterns: Visualization as scaffolding for ethnographic insight. Blog post for Ethnography Matters, <http://ethnographymatters.net/blog/2013/04/11/visualizing-plant-wars-player-patterns-to-aid-ethnography/> (last accessed December 13, 2021).
- Turing IBA (1950) Computing machinery and intelligence. *Mind; A Quarterly Review of Psychology and Philosophy* 59(236): 433.
- Wang T (2013) Big data needs thick data. *Ethnography Matters* 13. <http://ethnographymatters.net/blog/2013/05/13/big-data-needs-thick-data/>, accessed December 1, 2021.
- Wieringa M (2020) What to account for when accounting for algorithms: A systematic literature review on algorithmic accountability. In: FAT* 2020 - Proceedings of the 2020 conference on fairness, accountability, and transparency, pp. 1–18. <https://doi.org/10.1145/3351095.3372833>.
- Wilf E (2013) Toward an anthropology of computer-mediated, algorithmic forms of sociality. *Current Anthropology* 54: 6.
- Williams K (2018) Engineering ethnography. In: *Ethnography for a Data-Saturated World*, pp. 82–102. Manchester: Manchester University Press.