

The Shape of Tweets to Come: Automating Language Play In Social Networks

Tony Veale

Abstract. Twitter has proven itself a rich and varied source of language data for linguistic analysis. For Twitter is more than a popular new platform for social interaction via language; in many ways Twitter constitutes a whole new genre of text, as users adapt to its limitations (140 character “*tweets*”) and its novel conventions (e.g. re-tweeting, hashtags). Language researchers can harvest Twitter data to study how users convey meaning with affect, and how they achieve stickiness and virality with the texts they compose. But Twitter presents an opportunity of another kind to the computationally-minded language researcher, a *generative* opportunity to study how algorithmic models might exploit linguistic hypotheses to compose novel and meaningful micro-texts of their own. This computational turn allows researchers to go beyond merely descriptive models of playful uses of language such as metaphor and irony. It allows researchers to test whether their models embody a sufficiently algorithmic understanding of a phenomenon to facilitate the construction of a fully-automated computational system, one that can generate its wholly novel examples that are deemed acceptable by humans. This chapter presents and evaluates one such system, a *Twitterbot* named *@MetaphorMagnet* that generates, expresses and shares its own playful insights on Twitter.

Keywords: Irony, Metaphor, Twitter, Verbal Humour, Twitterbots

1. Introduction

A mismatch between a container and its contents can often tell us much more than the content itself, as when a person places the ashes of a deceased relative in a coffee can, or sends a brutal death threat in a *Hallmark* greeting card. The communicative effectiveness of mismatched containers is just one more reason to be skeptical of the *Conduit* metaphor (Reddy, 1979) – which views linguistic constructs as containers of propositional content to be faithfully shuttled between speaker and hearer – as a realistic model of human communication. Language involves more

than the faithful transmission of logical propositions between information-hungry agents, and more effective communication – of attitude, expectation and creative intent – can often be achieved by abusing our linguistic containers of meaning than by treating them with the sincerity that the *Conduit* metaphor assumes. Consider the case of verbal irony, in which a speaker deliberately chooses containers that are pragmatically ill-suited to the conveyance of their contents. For instance, the advertising container “*If you only see one [X] this year, make it this one*” assumes that [X] denotes a category of event – such as “*romantic comedy*” or “*movie about superheroes*” – with a surfeit of available members for a listener to choose from. When [X] is bound to the phrase “*comedy about Anne Frank*” or “*musical about Nazis*”, the container proves too hollow for its content, and the reader is signaled to the presence of irony. Though such a film may well be one-of-a-kind, the ill-fitting container suggests there are good reasons for this singularity that do not speak to X’s quality as an artistic event. Yet if carefully chosen, an apparently inappropriate container can communicate a great deal about a speaker’s relationship to the content conveyed within, and as much again about the speaker’s relationship to their audience.

As more practical limitations are placed on the form of linguistic containers, the more incentive one has to exploit or abuse containers for creative ends. Consider the use of Twitter as a communicative medium: writers are limited to micro-texts of no more than 140 characters to convey both their meaning and their attitude to this meaning. So each micro-text, or *tweet*, becomes more than a container of propositional content: each is a brick in a larger edifice that comprises the writer’s online personae and textual aesthetic. Many Twitter users employ irony and metaphor to build this aesthetic and thus build up a loyal audience of followers for their world view. Yet Twitter challenges many of our assumptions about irony and metaphor. Such devices must be carefully modulated if an audience is to perceive a speaker’s meaning in the playful (mis)match of a linguistic container to its contents. Failure to do so can have serious repercussions when one is communicating to thousands of followers at once, with tweets that demand concision and leave little room for nuance. It is thus not unusual for even creative tweets to come packaged with an explicit tag such as *#irony*, *#sarcasm* or *#metaphor*.

Metaphor and irony are much-analyzed phenomena in social media, but this chapter takes a generative approach, to consider the *production* rather than the *analysis* of creative linguistic phenomena in the context of a fully-autonomous computational agent – a *Twitterbot* – that crafts its own metaphorical and ironical tweets from its own knowledge-base of common-sense facts and beliefs. How might such a system exhibit a sense of irony that human users will find worthy of attention, and how might this system craft interesting metaphoric insights from a knowledge-base of everyday facts that are as banal as they are uncontentious? We shall explore the

variety of linguistic containers at the disposal of this agent – a fully autonomous computational system on the Web named *@MetaphorMagnet* – to better understand how such containers can be playfully exploited to convey ironic, witty or thought-provoking views on the world. With *@MetaphorMagnet* we aim to show that interesting messages are not crafted from interesting contents, or at least not necessarily so. Rather, effective tweets emerge from an appropriate if non-obvious combination of familiar containers with unsurprising factual fillers. In support of this view, we present an empirical analysis of the assessment of *@MetaphorMagnet*'s uncurated figurative outputs by human judges.

Just as one can often guess the contents of a physical container by its shape, one can often guess the meaning of a linguistic container by its form. We become habituated to familiar containers, and just as we might imagine our own uses for a physical container, we often pour our own meanings into suggestive textual forms. For in language, meaning follows form, and readers will generously infer the presence of meaning in texts that are well-formed and seemingly the product of an intelligent entity, even if this entity is *not* intelligent and any meaning is *not* intentional. Remarkably, Twitter shows that we willingly extend this generosity of interpretation to the outputs of bots that we *know* to be unthinking users of wholly aleatoric methods. Twitterbots exploit this *charity of interpretation*, wherein a well-formed linguistic container is assumed to carry a well-founded meaning – by serving up linguistic forms that readers tacitly fill with their own meanings. We aim to empirically demonstrate here that readers do more than willingly suspend their disbelief, and that a well-packaged linguistic form can seduce readers into seeing what is not there: a comprehensible meaning, or at least an intent to be meaningful. We do this by evaluating two metaphor bots side-by-side: a rational, knowledge-based Twitterbot named *@MetaphorMagnet* vs. an aleatoric and largely knowledge-free bot named *@MetaphorMinute*.

2. Digital Surrealists: Out of the Mouths of Bots

Most Twitterbots are simple, rule-based systems that use stochastic tools to explore a loosely-defined space of possible textual forms, or what Oulipo (1981) calls a space of “potential literature.” Most bots are thus high-concept, low-complexity generative systems that transplant the aleatoric methods and constraints of the early surrealists, the Oulipo group and the “beat” writers – from André Breton to Raymond Queneau (1961, 1981) to William Burroughs (1963) and Brion Gysin – into the realms of digital content, social networks and online publishing. Each embodies a language game with its own generative rules, or what Breton called “*la règle du jeu*.” Yet Breton, Queneau, Burroughs and Gysin viewed the use of mechanical

rules as merely the first of a two-stage text-creation process: in the first, random recombinant methods are used to confect candidate texts in ways that, though unguided by meaning, are also free of the baleful influence of cliché; at the second stage, these candidates are carefully interpreted and filtered by a human, to select those that are novel and interesting and to reject the rest. Most bots implement the first stage but ignore the second, pushing the task of critiquing and filtering candidate texts onto the humans who read and selectively re-tweet them.

Nonetheless, some bots achieve surprising linguistic results with the simplest of tools. Consider *@Pentametrone*, a bot that generates accidental poetry by re-tweeting pairs of random tweets of ten syllables apiece (for an iambic pentameter reading) if each ends on a rhyming syllable. When the meaning of each tweet in a couplet happens to cohere with the other, as in “*Pathetic people are everywhere*”/“*Your web-site sucks, @RyanAir*”, the pairing produces an emergent meaning that is richer and more resonant than that of either tweet alone. Trending social events such as the *Oscars* or the *Super Bowl* are especially conducive to just this kind of synchronicity, as in this fortuitous pairing: “*So far the @SuperBowl commercials blow*” / “*Not even gonna watch the halftime show.*” If one were so minded, one could trace the lineage of *@Pentametrone* and that of other tweet-splicing bots (such as the headline-creating *@twoheadlines* of Darius Kazemi) to Oulipo mechanisms such as Queneau’s (1961) flip-book for generating sonnets.

In contrast, a bot named *@MetaphorMinute* wears its aleatoric methods on its sleeve, for its tweets – such as “*a haiku is a tonsil: peachblow yet snail-paced*” – are not so much random metaphors as random metaphor-shaped texts. Much like Chamberlain and Etter’s (1983) RACTER system, *@MetaphorMinute* employs a generative strategy that relies on word associations and randomness rather than word meanings and pragmatics. The bot instantiates a standard linguistic container for metaphors – the copula frame “*X is a Y*” – with random-seeming word choices, and tweets the results without any attempt at quality estimation or filtering every two minutes. Interestingly, *@MetaphorMinute*’s tweets are just as likely to provoke a sense of mystification and ersatz profundity as they are total incomprehension. Bots such as *@Pentametrone* and *@MetaphorMinute* do not generate their texts from the semantic-level up; rather, they manipulate texts at the word-level only, and thus lack any sense of the meaning of a tweet, or any rationale as to why one tweet might be better – which is to say, more interesting, more apt or more re-tweetable – than any other.

The Full-FACE poetry generator of Colton *et al.* (2012) also uses a template-guided version of the cut-up method to mash together semantically-coherent text fragments in a way that – much like *@Pentametrone* – obeys certain over-arching constraints on metre and rhyme. These text fragments come from a variety of online sources, ranging from short tweets to long news articles. News stories are a rich source of

readymade phrases that convey resonant images, and these can be clipped from a news text using standard NLP techniques, while tweets that use affect-rich language can also be extracted automatically via standard *sentiment analysis* lexica and tools. Thus, a large stock of resonant similes, such as “*blue as a blueberry*” or “*hot as a sauna*” can be extracted from the Web using a search engine (Veale, 2014), since the simile frame “*as X as Y*” is specific enough to query for, and promiscuous enough to match, a rich diversity of typical X:Y associations. These associations can then be recast in a variety of poetic forms to make their clichéd offerings seem fresh again, as in “*Blueberry-blue overalls*” or “*sauna-hot jungle*.”

Indeed, the very act of juxtaposing clichés can itself be a creative act, as evidenced both by the success of the cut-up method in general and that of specific cut-ups in particular. Consider William Empson’s withering analysis of the persnickety, cliché-hating George Orwell, whom Empson called “*the eagle eye with the flat feet*” (quoted in Ricks [1995:356], who admires Empson’s “*audacious compacting of clichés*”). The Full-FACE system is just one of many computational creativity systems that use an autonomous variant of Burroughs and Gysin’s cut-up method to integrate tight constraints on form with loose constraints on meaning.

Breton famously stated that “*Je ne veux pas changer la règle du jeu, je veux changer de jeu*.” Twitterbots do not change or transcend their own rules, but different bots do represent different language games with their own rules. So to change the game, a computational creativity developer can simply build a new bot, to exploit a different set of tropes and linguistic containers. It is rare for any one Twitterbot to incorporate a diverse set of tropes and production mechanisms; each typically follows Breton’s experimentalist approach to art in its random sampling of a specific space of possibilities. Each bot thus forms its own art installation, to showcase a single generative idea. @*MetaphorMagnet*, the bot at the heart of this chapter, represents a departure from this norm, insofar as it exploits a wide range of tropes and rendering strategies, it employs diverse sources of knowledge, and it applies a variety of reasoning styles to generate surprising conclusions from its stock of otherwise banal facts. But does this added sophistication – bought at the cost of increased system complexity and knowledge-engineering effort – result in tweets that are seen as more meaningful, novel, apt or retweetable by human users? It is this point that exercises us most in the coming sections.

3. Exercises in Style Over Substance

Style lends a distinctive shape and appearance to our linguistic containers. But it can do more than make one text stand out from others: it can shape the way one should feel about the contents of a linguistic container. Style can imbue a banal event with the excitement of a thriller, the immediacy of

a newsflash or the comedy of a farce. It can make a crass generalization read like a dry scientific fact, or a dull fact read like a tabloid headline. As shown by Raymond Queneau's remarkable *Exercises in Style* (1947/1981), style crucially shapes how readers construe the states of affairs conveyed by a text. Queneau at turns wrings pathos, humour and cultural insight from his 99 alternate stylistic renderings of the same banal tale of a scuffle on a bus. Though most Twitterbots embody just a single exercise in style, these bots collectively offer a digital realization of Queneau's entire experimental agenda. Indeed, it is not too great a stretch to suggest that humans follow these bots for many of the same reasons they read Queneau's *Exercises* or his other Oulipo work (such as *Cent mille milliards de poèmes* from 1961). Meanwhile, more ambitious bots are taking up the Queneau challenge to create original outputs in a wide diversity of styles that grab the eye and shape the attitude of a reader. @BestOfBotWorlds, for instance, invents tweets with faux-inspirational or religious content using forms that satirize the "style" of Jesus, Mohammad, Yoda, Donald Trump and even the Hulk.

We humans obtain more mileage than we may ever care to admit from templates, constraints and other "bot"-like stylistic approaches to linguistic creativity. Consider what Matthew McGlone and Jessica Tofighbakhsh (1999) call the *Keats heuristic*, an insight into creative language use that owes as much to Nietzsche ("*we sometimes consider an idea truer simply because it has a metrical form and presents itself with a divine skip and jump*") as to the poet John Keats ("*Beauty is truth, truth beauty*"). McGlone and Tofighbakhsh (2000) show that when presented with uncommon maxims or proverbs with internal rhyme (e.g. "*woes unite foes*"), subjects tend to view these as more insightful about the world than the equivalent paraphrases with no internal rhyme at all (e.g. "*troubles unite enemies*"). While the Keats heuristic is not exactly a license to pun, it is an incentive to rhyme, and to give as much weight (or more still) to superficial aspects of poetry generation as to deep semantics and pragmatics. Indeed, the heuristic is tacitly central to the operation of virtually every computational creativity approach to poetry generation (e.g. Milic, 1970; Chamberlain and Etter, 1983; Gervás, 2000; Manurung *et al.* 2012; Veale, 2013). If human poets ask questions first and rhyme later, computational creativity systems typically rhyme first and ask questions later, if at all. For if the human jury in the O.J. Simpson trial could be turned against bald facts with the Keatsian "*If the glove don't fit you must acquit*", readers of computer-generated poetry can be persuaded to see deliberate meaning and resonance in any output that has a "*divine skip and jump*."

There is something undeniably special about poetry, whether it is the gentle poetry of William Shakespeare's "*Shall I compare thee to a summer's day*" or the rough poetry of Johnnie Cochrane's "*If the glove don't fit you must acquit*". Milic (1970), an early computational creativity pioneer, argues that while poetry "*is more difficult to write than prose*" it

offers other freedoms to writers due to the willingness of readers to “interpret a poem, no matter how obscure, until he has achieved a satisfactory understanding.” What then of the enigmatic tweets of bots like @MetaphorMinute, whose obscurity is a function of random word choice and whose surface forms are not designed to make any sense at all? Milic argues that computer poetry serves a useful role other than its obviously generative one, by alerting us to “the curious behavior of familiar words in unfamiliar combinations.” Behaviour that makes perfect sense when dealing with the writings of a gifted human poet, such as our tendency to “interpret an utterance by making what concessions are necessary on the assumption that a writer has something in mind of which the utterance is the sign”, is, argues Milic, “inappropriate when the speaker is a computer.” Yet Twitterbots benefit from such concessions and assumptions whether or not followers know them to be bots. This *Eliza effect* (see Weizenbaum, 1966; Hofstadter, 1995) is especially pronounced in the coining of would-be metaphors, leading Milic to note “how readily we accept metaphor as an alternative to calling a sentence nonsensical.” @MetaphorMinute and other aleatoric bots wring maximal value from this insight by devising texts that they themselves cannot distinguish from nonsense. This begs an important question: are the meanings imposed on a random text by a creative human of comparable value to those conveyed by a Twitterbot with its own model of the world and its own insights to tweet?

4. Filling Linguistic Containers With Metaphorical Meanings

What might it mean for a bot to have “something in mind of which [its] utterance is the sign”? When it comes to metaphor generation, we might expect that our bot would generate its figurative tweets from a conceptual model of the world as it sees it, in a way that accords with a sound theory of *how* and *why* humans actually use metaphor. For the latter, the field of Artificial Intelligence offers us a range of models to choose from.

Computational approaches to metaphor divide into four broad classes: the categorial, the corrective, the analogical and the schematic. Categorial approaches view metaphor as a means to re-conceptualize one idea by placing it into a taxonomic category strongly associated with another (see Hutton, 1982; Way, 1991; Glucksberg, 1998). Corrective approaches view metaphor as an inherently anomalous deviation from literal language, and strive to *recover* the corresponding literal meaning of any figurative statement that violates its lexico-semantic norms (see Wilks, 1978; Fass, 1991). The analogical approaches aim to capture the relational parallels that allow our representation of an idea in one domain, the *source*, to be systematically projected onto our mental representation of an idea in another, the *target* (see Gentner *et al.*, 1989; Veale and Keane, 1997).

Finally, schematic approaches aim to explain how related linguistic metaphors arise as surface manifestations of deep seated cognitive structures called *Conceptual Metaphors* (Lakoff and Johnson, 1980; Carbonell, 1981; Martin, 1990; Veale and Keane, 1992). Each approach has its own merits, but none offers a complete computational solution. Bots that aim for a general competence in metaphor must thus implement a selective hybrid of multiple approaches. Yet each approach also requires its own source of knowledge. Categorical approaches require a comprehensive taxonomy of flexible categories that can embrace atypical members on demand. Corrective approaches are built on a substrate of literal case-frames onto which deviant usages can be correctively projected. Analogical approaches assume an inventory of graph-theoretic representations of concepts, from which a structure-mapping engine can eke out its sub-graph isomorphisms. Schematic approaches rely on a stock of Conceptual Metaphors (CMs) – such as *Life is a Journey* or *Theories are Buildings* – to unearth the deep structures beneath the surface of diverse linguistic forms.

Though hybrid approaches demand multiple sources of knowledge, there exist public Web services that integrate this knowledge with the appropriate means of using it for metaphor. The *Thesaurus Rex* Web service of Veale and Li (2013) provides a highly divergent system of fine-grained categorizations that allows a 3rd-party client system to e.g. determine that *War* and *Divorce* have each been viewed as kinds of *destructive thing*, *traumatic event* and *severe conflict* in the texts of the Web. The *Metaphor Eyes* Web service of Veale and Li (2011) is a rich source of relational norms – also harvested at scale from Web texts – such as that businesses earn profits and pay taxes, or that religions ban alcohol and believe in reincarnation. The *Metaphor Magnet* service of Veale (2014) offers a rich source of the stereotypical properties and behaviors of familiar ideas, and provides a means to retrieve salient CMs from the Google n-grams (Brants and Franz, 2006) which can then be elaborated to create novel linguistic metaphors.

@MetaphorMagnet relies on each of these public Web services to generate the conceptual conceits that underpin its figurative tweets. For instance, it uses *Thesaurus Rex* to provide the categorization insights that it then packages as *odd-one-out* lists or as *faux-dictionary* definitions. It uses the *Metaphor Eyes* service to provide the relational structures it needs to perform structure mapping and thus concoct original analogies and dis-analogies. And it uses the *Metaphor Magnet* service to access the stereotypical properties and behaviors of ideas, and to juxtapose these properties via resonant contrasts and norm contraventions. Once the conceptual chassis of a metaphor is constructed in this way, it is then packaged in an apt linguistic form.

5. Shaping a Tweet: Automated Exercises in Twitter Style

CMs such as *Life Is A Journey* and *Politics Is A Game* are more than productive deep-structures for the generation of whole families of linguistic metaphors; they also provide the conceptual mappings that shape our habitual thinking about such familiar ideas as *Life, Love, Politics* and *War*. Politicians and philosophers exploit conceptual metaphors to frame an issue and shape our expectations; when a CM fails to match our own experience, we reject it and switch to a more apt metaphor. So a metaphor-generating bot can thus create a thought-provoking opposition by pitting one CM against another that advocates a conflicting view of the world. The following tweet from *@MetaphorMagnet* uses this approach to contrast two views on *#Democracy*:

*To some voters, democracy is an important cornerstone.
To others, it is a worthless failure.*
#Democracy=#Cornerstone #Democracy=#Failure

The CM *Democracy Is A Cornerstone* (of society) is often used to frame political discussions, and can be seen as an specialization of the CM *Society Is A Building*, itself an elaboration of the CM *Organization Is Physical Structure* (see Grady, 1997). Yet the importance of cornerstones to the buildings they anchor finds a sharp contrast in the assertion that *Democracy Is A Failure*. Each of these affective claims is so commonly asserted that they can be found in the Google n-grams, a large database of short fragments of frequent Web texts. The 4-gram “*democracy is a cornerstone*” has a frequency of 91 in the Google n-grams, while the 4-gram “*democracy is a failure*” has a frequency of 165. These n-grams, which suggest potential CMs for *@MetaphorMagnet*, are elaborated with added detail via the *Metaphor Magnet* Web service, which tells the bot that the stereotypical *cornerstone* is *important* and the stereotypical *failure* is *worthless*. The following tweet makes similar use of CMs found in the Google n-grams, but renders the conflict in a different linguistic container:

*Remember when tolerance was promoted by crusading liberals?
Now, tolerance is violence that only fearful appeasers can avoid.*

The bot is guided here by the suggestive Google 3-gram “*Tolerance for Violence*” (frequency=1353), but it does not directly contrast the ideas *#Tolerance* and *#Violence*. Instead, it finds a potential analogy in this juxtaposition, between the promoters of *#Tolerance* (which it renders as *crusading liberals*) and the opponents of *#Violence* (which it renders as *fearful appeasers*). The choice of stereotypical properties (*crusading* and *fearful*) is driven by the bot’s need to create a resonant semantic opposition.

The bot omits the hashtags [#Tolerance](#)=[#Violence](#) from this tweet due to the confines of Twitter's 140-character limit. But it can also choose to render a complex conceit across two successive tweets, as in the following:

Remember when research was conducted by prestigious philosophers?
[#Research](#)=[#Fruit](#) [#Philosopher](#)=[#Insect](#)

Now, research is a fruit eaten only by lowly insects.
[#Research](#)=[#Fruit](#) [#Philosopher](#)=[#Insect](#)

[@MetaphorMagnet](#) uses a number of packaging strategies to turn a figurative comparison into an ironic observation, ranging from the use of an explicit [#Irony](#) hashtag (which is commonplace on Twitter) to the use of “scare” quotes to focus on the part of a tweet most deserving of disbelief. The following tweet showcases both of these strategies:

[#Irony](#): *When some chefs prepare "fresh" salads the way apothecaries prepare noxious poisons.*
[#Chef](#)=[#Apothecary](#) [#Salad](#)=[#Poison](#)

Irony offers a concise means of contrasting two points of view: that which is expected and the disappointing reality. By comparing the preparation of salads – the “*healthy*” option on most menus – to the preparation of poisons, this analogy undermines the expectation of healthfulness and suggests that some salads are noxious and chemical-filled. The real world is filled with situations in which naturally antagonistic properties are found in surprising proximity. These situations, if expressed in the right linguistic form, can be elevated to the level of situational irony. Consider, for instance, the following [@MetaphorMagnet](#) tweet:

[#Irony](#): *When the timers that are found in enjoyable games activate gruesome bombs.* [#Enjoyable](#)=[#Gruesome](#)

It is important to stress that [@MetaphorMagnet](#) does not simply fill linguistic templates with related words. Rather, the above tweet is constructed at the knowledge-level, by a bot that intentionally seeks out stereotypical norms that are related (e.g. by a pivotal idea *timer*) yet which can be placed into antagonistic juxtapositions around this pivot. In effect, the goal of the linguistic rendering is to package a knowledge-level conceit – typically a conflict of ideas and properties – in a tweet-sized narrative. For example, the following tweet is stylistically rendered as a narrative of change:

To join and travel in a pack: This can turn pretty girls into ugly coyotes.
[#Girl](#)=[#Coyote](#)

Twitter offers unique social affordances that allow a bot to elevate almost any contrast of ideas into a dramatic narrative. Rather than talk of generic liberals or appeasers, a bot can give these straw men real names, or at least invent fake names that look like the real thing and which, as Twitter handles, seem wittily apropos to the views that are espoused. In this way, by imagining its central conceit as a topic of a vigorous debate by real people, a bot can turn an abstract metaphor into a concrete situation with its own colorful participants. Consider the social debate that is made personal in this tweet from [@MetaphorMagnet](#):

[.@war_poet](#) says history is a straight line
[.@war_prisoner](#) says it is a coiled chain
[#History=#Line](#) [#History=#Chain](#)

The handles [@war_poet](#) and [@war_prisoner](#) are invented by [@MetaphorMagnet](#) to suit, and amplify, the figurative views that they are advanced in the tweet, by using a mix of relational knowledge (from the *Metaphor Eyes* service) and language data (via the Google n-grams). Since poets write poems about the wars that punctuate history, and poems contain lines, the 2-gram “*war poet*” is recognized as an apt handle for an imaginary Twitter user who might advance a view of *history as a line*. In this case the handle [@war_poet](#) really does name a real Twitter user, but this only adds to the sense that Twitterbot confections are a new kind of interactive theatre and performance art (see Dewey, 2014). The most profound aspects of this contrast are not appreciated by [@MetaphorMagnet](#) itself, or at least not yet. For example, the bot does not yet appreciate what it means for history to be a straight line, and while it knows enough to invent the intriguing handle [@war_prisoner](#), neither does it appreciate what it might mean to be a prisoner of history, enslaved in a repeating cycle of war. Our bots will always evoke in we humans more than they themselves can ever appreciate, yet this may itself be a key part of computational creativity’s allure.

6. Content Versus Container: Evaluating Metaphor Generation

[@MetaphorMagnet](#) differs from [@MetaphorMinute](#) in a number of key ways. For one, its mechanics are informed by Lakoff and Johnson’s *Conceptual Metaphor Theory* and a range of computational approaches. For another, it draws on considerable semantic and linguistic resources, from a large knowledge-base of conceptual relations and stereotypical beliefs to the linguistic diversity of the Google n-grams. Note that *all* of

[@MetaphorMagnet](#)'s tweets – all its hits and all its misses – are open to public scrutiny on Twitter. But to empirically evaluate the success of the bot as a knowledge-based, theory-driven producer of novel, meaningful and retweet-worthy metaphors, we turn to the crowdsourcing platform [CrowdFlower](#), where we conduct a comparative evaluation of [@MetaphorMagnet](#) and its closest knowledge-free counterpart, [@MetaphorMinute](#). The latter, designed by noted bot-maker Darius Kazemi, uses a wholly aleatoric approach to metaphor generation yet has over 500 followers on Twitter that do not mind its *one-every-two-minutes* scattergun approach to generation. [@MetaphorMinute](#) crafts metaphors by filling a template with nouns and adjectives that are chosen more-or-less at random, to produce inscrutable tweets such as “*a cubit is a headboard: stational yet tongue-obsessed.*”

We chose 60 tweets at random from the past outputs of each Twitterbot. CrowdFlower annotators, who were each paid a small sum per judgment, were not informed of the origin of any tweet, but simply told that each was selected from Twitter because of its metaphorical content. We did not want annotators to actively suspend their disbelief by knowingly dealing with bot outputs. Annotators were paid to rate the content of each tweet along three dimensions, *Comprehensibility*, *Novelty* and likely *Retweetability*, and to rate all three dimensions on the same scale: *Very Low* to *Medium Low* to *Medium High* to *Very High*. Ten annotations were solicited for each dimension of each tweet, though the responses of likely scammers (non-engaged annotators) were later removed from the dataset. Tables 1 through 3 present the distributions of mean ratings per tweet, for each dimension and each Twitterbot..

<i>Comprehensibility</i>	<i>Metaphor Magnet</i>	<i>Metaphor Minute</i>
<i>Very Low</i>	11.6%	23.9%
<i>Med. Low</i>	13.2%	22.2%
<i>Med High</i>	23.7%	22.4%
<i>Very High</i>	51.5%	31.6%

Table 1. Comprehensibility of [@MetaphorMagnet](#) and [@MetaphorMinute](#)

More than half of [@MetaphorMagnet](#)'s tweets were ranked as having very high comprehensibility, while less than one third of [@MetaphorMinute](#)'s

tweets are so ranked. More surprising, perhaps, is the result that annotators found more than half of [@MetaphorMinute](#)'s wholly random metaphors to have medium-high to very-high comprehensibility. This Twitterbot's use of abstruse terminology, such as *stational* and *peachblow*, may be a factor here, as might the bot's use of the familiar copula container *X is Y* for its metaphors, which may well seduce annotators into believing that an apparent metaphor really does have a comprehensible meaning, if only one were to expend enough mental energy to actually discern it.

<i>Novelty</i>	<i>Metaphor Magnet</i>	<i>Metaphor Minute</i>
<i>Very Low</i>	11.9%	9.5%
<i>Med. Low</i>	17.3%	12.4%
<i>Med High</i>	21%	14.9%
<i>Very High</i>	49.8%	63.2%

Table 2. Novelty ratings of [@MetaphorMagnet](#) and [@MetaphorMinute](#)

The dimension *Novelty* yields results that are equally surprising. While half of [@MetaphorMagnet](#)'s metaphors are rated as having very-high novelty in Table 2, almost two-thirds of [@MetaphorMinute](#)'s tweets are just as highly rated. However, we should not be overly surprised that [@MetaphorMinute](#)'s bizarre juxtapositions of rare or unusual words, as yielded by its unconstrained use of aleatoric techniques, are seen as more unusual than those word juxtapositions arising from [@MetaphorMagnet](#)'s controlled use of attested Web n-grams and stereotypical knowledge. As shown by Giora *et al.* (2004), novelty is neither a source of pleasure in itself nor is it a reliable benchmark of creativity. Rather, pleasurability derives from the recognition of *useful* novelty, that is, novelty that can be understood and appreciated relative to the familiar.

On Twitter, useful exploitation is frequently a matter of social reach. A tweet is novel and useful to the extent that it attracts the attention of Twitter users and is deemed worthy of re-tweeting to others in one's social circle. Our third dimension, *Re-Tweetability*, reflects the likelihood that an annotator would ever consider re-tweeting a given metaphorical tweet to others. Though we ask annotators to speculate here – neither bot has enough followers to perform a robust statistical analysis of actual retweet

rates – the results largely conform to our expectations. The results presented in Table 3 show retweetability to be a matter of novelty *and* comprehensibility together, and not just a matter of novelty alone. Though annotators are not generous with their *Very-High* ratings for either bot, [@MetaphorMagnet](#)'s tweets are judged to be significantly more retweetable than the largely random offerings of [@MetaphorMinute](#).

<i>Retweetability</i>	<i>Metaphor Magnet</i>	<i>Metaphor Minute</i>
<i>Very Low</i>	15.5%	41%
<i>Med. Low</i>	41.9%	34.1%
<i>Med High</i>	27.4%	15%
<i>Very High</i>	15.3%	9.9%

Table 3. Retweetability of [@MetaphorMagnet](#) and [@MetaphorMinute](#)

This is just as well, given the considerable gap in complexity and sophistication that exists between the two bots. But this is an encouraging result not just for theory-informed Twitterbots like [@MetaphorMagnet](#) and their creators, but for Twitter itself. Twitter offers a compelling platform for research into interactive play through language, not least because its human users appreciate these phenomena when they see them.

But as pointed out in Reddy (1979), the conduit metaphor of language is an imperfect one, and the linguistic containers that shuttle back and forth between speakers may convey much more or much less than they appear to convey. The appearance of comprehensibility may not always result in actual comprehension, and so, while a Computational Creativity system may cleverly use packaging and style to foster a belief that a given tweet has a coherent meaning, it cannot insert this meaning into the head of a reader. Meaning is the product of interpretation, and interpretation is often hard. Milic (1970) notes that in a context that licenses a poetic interpretation, such as one in which a reader is told that a particular text is a metaphor, readers are more likely to accept that the text – as inscrutable as it may be – has a metaphorical meaning rather than dismiss it as nonsense. Recall that over 75% of [@MetaphorMagnet](#)'s tweets and over 50% of [@MetaphorMinute](#)'s tweets are judged as having *medium-high* to *very-high*

comprehensibility. We thus need to look deeper, to determine whether raters can actually back up these judgments with actual meanings.

So in a second CrowdFlower experiment, we make raters work harder, to reconstruct a partial tweet by adding the missing information that will make it whole and apt again. That is, we employ a *cloze* test format for this experiment, by removing from each tweet the pair of key qualities that anchor the tweet and make its comparison of ideas seem meaningful and apt. For *@MetaphorMagnet*, for example, we remove the properties *detailed* and *vague* in this tweet:

*To some freedom fighters, freedom is a detailed recipe.
To others, it is a vague dream.*
[#Freedom=#Recipe](#) [#Freedom=#Dream](#)

For *@MetaphorMinute*, we blank the qualities *hippy* and *revisional* in the following tweet:

a flatfoot is a houseboat: hippy and revisional

For each tweet from each bot, we blank out a pair of original qualities as above; this pairing is the answer that is sought from human judges. We also choose 4 distractor pairs for each original pair, by selecting pairs from other tweets from the same bot. As in our first experiment, we chose 60 tweets at random from the past outputs of each bot, and 10 ratings were solicited for each. Annotators were presented with a tweet in which the key properties were blanked out, and given five randomly ordered pairs of possible fillers (the original pair and four distractors from other tweets) to choose from.

<i>Aptness</i>	<i>Metaphor Magnet</i>	<i>Metaphor Minute</i>
<i>Very Low</i>	0%	84%
<i>Med. Low</i>	22%	16%
<i>Med High</i>	58%	0%
<i>Very High</i>	20%	0%

Table 4. Relative Aptness of [@MetaphorMagnet](#) and [@MetaphorMinute](#)

To make the results of the experiment comparable to those of the 1st experiment (Tables 1,2,3), we obtain the mean aptness of each tweet, so that e.g. if 7 out of 10 raters correctly choose the original pairing, then that tweet is deemed to have an aptness of 0.7. We then place these aptness scores into bands, where the *Very Low* band = 0 to 0.25, *Medium Low* = 0.26 to 0.5, *Medium High* = 0.51 to 0.75, and *Very High* = 0.76 to 1. By calculating the distribution of tweets to each band, we can determine e.g. the percentage of tweets from each bot that are put into the *Very High* band.

Our hypothesis is rather straightforward: if tweets are linguistic containers that are carefully crafted to convey a particular meaning, then it should be easier to select the missing pair of qualities that make this meaning whole again; if, on the other hand, the tweet is all there is, and its content is chosen mostly at random, then raters will choose the right pairing with no more success than random selection. The results reported in Table 4 bear out this hypothesis.

The *Eliza* effect (Hofstadter, 1995) can lead us to appreciate a bot's tweets as meaningful but it cannot tell us what this meaning should be. Though the results above may seem a foregone conclusion, as *@MetaphorMagnet's* tweets are designed to communicate a fully recoverable meaning while those of *@MetaphorMinute* are not, this is surely what it means to engage in real communication: to design an utterance so that an intended meaning is re-created, in whole or in part, in the mind of an intelligent, receptive audience.

7. Concluding Remarks: Believing is Seeing

Whenever a machine uses style and packaging to convey a sense of understanding and profundity with otherwise shallow linguistic forms, as in Weizenbaum's (1965) infamous ELIZA system (which fooled some users into believing it was a fully-functional and caring psychotherapist), the label "*ELIZA Effect*" proves to be an apt one (Hofstadter, 1995). However, we humans are also subject to an *ELIZA effect* of our own, insofar as we often do others the courtesy of assuming their utterances to be freighted with real meaning and creative intent, and will often work hard to uncover that meaning for them.

At one time or another, we have all relied on catch-phrases, clichés, slogans, idioms, canned jokes and other half-empty linguistic containers to suggest to others that we have deeper meanings in mind, or have something more profound to offer, than we actually do. In a famous polemical essay from 1946, George Orwell excoriates speakers of English for their reliance on jargon, foreign words and empty phraseology as a substitute for thoughts of real substance, while Geoff Pullum (2003) upbraids modern

speakers for a grating over-reliance on “*multi-use, customizable, instantly recognizable, time-worn, quoted or misquoted phrases or sentences that can be used in an entirely open array of different jokey variants by lazy journalists and writers.*” These “*phrases for lazy writers in kit form*” are not that different from the template-based language games played by superficial Twitterbots, and though we humans fill our templates – such as “*X is the new black*”, “*In X no one can hear you scream*” or “*if the Eskimos have N words for snow then Xs surely have as many for Y*” – with lexical fillers that are contextually apt, we employ our templates to be just as provocative, and to imply or to suggest more than we actually mean. By fitting our meanings to familiar structures, we give our readers the cues they need to see the meanings we want them to see.

As the cliché goes, *seeing is believing*, yet when metaphor and irony are involved, the converse is also true, for these turn believing into a kind of *seeing*. Ultimately, metaphor and irony are cognitive devices for generating and conveying a specific viewpoint. So at the level of ideas, every conceptual metaphor offers a means of viewing one idea through the lens of another, while at the level of words every linguistic metaphor offers a means of concisely conveying this viewpoint to another, and of helping others to view the world from the same vantage point. If the situation with irony is that much more complicated, it is because irony offers a *stereoscopic* viewpoint that conveys and contrasts *two* conflicting perspectives at once to highlight the disparity between the world as it appears and the world as it should be, or at least as it was promised to be. Yet whether one is thinking and communicating via metaphor or via irony, construction proceeds in much the same direction, from perception to conceptualization to expression.

The hashtag *#irony* is more often used on Twitter as an observation than as a warning, and corresponds more to the conversational gambit “*Isn't it ironic ...*” than to an insurance policy against potential misunderstanding. Yet a creative writer cannot make a tweet ironic by adding the hashtag *#irony* any more than one can make it a metaphor by adding *#metaphor* or turn it into a witticism by adding *#funny*. Though first-generation bots such as *@MetaphorMinute* achieve the appearance of metaphor, by generating random *metaphor-shaped* tweets that merely use the copula container for its *X is Y* metaphors, there are no containers that can imbue random juxtapositions with ironic or figurative meanings, and no hashtag that can magically substitute for a lack of original insight. To generate a creative figurative statement that really is designed to be playfully figurative, such as a witticism that is both ironic and metaphorical, an automated system must proceed in the same way as a human: from insight to conceptualization to careful packaging in an appropriate linguistic container. In the chapter we have shown how an automated system that is both theory-guided and knowledge-driven – a Twitterbot called

@MetaphorMagnet – navigates this course for itself. Modular theoretical frameworks, such as the *General Theory of Verbal Humour* (GTVH) used here, have their sectional boundaries severely tested by the task of squeezing provocative figurative forms into the textual confines of a tweet. When packing so much into so little, every single aspect of language production – from scenario construction to schema selection to framing strategy to word choice – must work so closely with every other that no single concern can ever be truly autonomous or “modular.” Twitter is not just the best of bot worlds then, but an ideal environment in which to study the cognitive linguistics of playful language.

8. References

- Thorsten Brants and Alex Franz.
 2006 Web 1T 5-gram database, Version 1. *Linguistic Data Consortium*.
- William S. Burroughs
 1963 The Cut-Up Method. LeRoi Jones (Ed.), *The Moderns: An Anthology of New Writing in America*. New York: Corinth.
- Jaime G. Carbonell
 1981 Metaphor: An inescapable phenomenon in natural language comprehension. *Report 2404*. Carnegie Mellon Computer Science Dept.
- William Chamberlain and Thomas Etter
 1983 *The Police-man’s Beard is Half-Constructed: Computer Prose and Poetry*. Warner Books.
- Simon Colton, Jacob Goodwin and Tony Veale.
 2012 Full-FACE Poetry Generation. In *Proc. of the 3rd International Conference on Computational Creativity*, Dublin, Ireland.
- Caitlin Dewey
 2014 What happens when @everyword ends? Intersect, *Washington Post*, [May 23rd edition](#).
- Dan Fass
 1991 Met*: a method for discriminating metonymy and metaphor by computer. *Computational Linguistics*, 17(1):49-90.
- Dedre Gentner, Brian Falkenhainer and Janice Skorstad
 1989 Metaphor: The Good, The Bad and the Ugly. In *Theoretical Issues in NLP*, Yorick Wilks (Ed.) Hillsdale, NJ: Lawrence Erlbaum Associates.
- Pablo Gervás.
 2000 Wasp: Evaluation of different strategies for automatic generation of Spanish verse. In *Proc. of the AISB-2000 Symposium on Creative & Cultural Aspects of AI*, 93-100.

- Rachel Giora, Ofer Fein, Jonathan Ganzi, Natalie A. Levi and Hadas Sabah
 2004 Weapons of Mass Distraction: Optimal Innovation and Pleasure Ratings. *Metaphor and Symbol* **19**(2):115-141.
- Sam Glucksberg
 1998 Understanding metaphors. *Current Directions in Psychological Science*, 7:39-43.
- Joseph Grady
 1997). Foundations of Meaning: Primary Metaphors and Primary Scenes. University of California.
- Douglas Hofstadter
 1995 The Ineradicable Eliza Effect and Its Dangers. *Fluid Concepts and Creative Analogies: Computer Models of the Fundamental Mechanisms of Thought* (Preface 4), Basic Books: New York.
- James Hutton (translator)
 1982 *Aristotle's Poetics*. New York, NY: Norton.
- George Lakoff and Mark Johnson
 1980 *Metaphors We Live By*. Chicago, Illinois: Chicago University Press.
- James H. Martin
 1990 *A Computational Model of Metaphor Interpretation*. Academic Press.
- Ruli Manurung, Graeme Ritchie and Henry Thompson
 2012 Using genetic algorithms to create meaningful poetic text. *JETAI* 24(1):43–64.
- Matthew S. McGlone and Jessica Tofighbakhsh
 1999 The Keats heuristic: Rhyme as reason in aphorism interpretation, *Poetics* **26**(4):235-44.
- Matthew S. McGlone and Jessica Tofighbakhsh.
 2000 Birds of a feather flock conjointly (?): rhyme as reason in aphorisms. *Psychological Science* **11** (5): 424–428.
- Louis T. Milic
 1971 The possible usefulness of computer poetry. *The Computer in Literary and Linguistic Research*, R.A. Wisbey (Ed.), Cambridge, MA.
- OULIPO
 1981 Atlas de littérature potentielle. Number vol. 1 in Collection Idées. Gallimard.
- George Orwell
 1946 Politics and the English language. *Horizon*, 13(76), April issue.
- Geoffrey Pullum.
 2003 Phrases For Lazy Writers in Kit Form. *Language Log post*, [October 27, 2003](#).

Raymond Queneau

- 1961 *Cent mille milliards de poèmes*. Schoenhof's Foreign Books, Inc.
 1981 *Exercises in Style. 2nd Edition* (Translated from the French by Barbara Wright). New York: New Directions Books.

Michael J. Reddy

- 1979 *The conduit metaphor: A case of frame conflict in our language about language*. In A. Ortony (Ed.), *Metaphor and Thought*, 284–310. Cambridge University Press.

Christopher B. Ricks

- 1980 Clichés. In: L. Michaels and C. Ricks (Eds), *The State of the Language*. University of California Press, Berkeley.

Tony Veale and Mark T. Keane

- 1992 Conceptual Scaffolding: A spatially founded meaning representation for metaphor comprehension. *Computational Intelligence* 8(3):494-519.

Tony Veale and Mark T. Keane

- 1997 The Competence of Sub-Optimal Structure Mapping on 'Hard' Analogies. In *Proceedings of IJCAI'97, the 15th International Joint Conference on Artificial Intelligence*. Nagoya, Japan. Morgan Kaufmann.

Tony Veale and Guofu Li.

- 2011 Creative Introspection and Knowledge Acquisition. In Proc. of AAAI-2011, *The 25th Conference of the Association for the Advancement of Artificial Intelligence*. San Francisco: AAAI Press.

Tony Veale and Guofu Li.

- 2013 Creating Similarity: Lateral Thinking for Vertical Similarity Judgments. In *Proceedings of ACL 2013, the 51st Annual Meeting of the Assoc. for Computational Linguistics, Sofia, Bulgaria*,

Tony Veale

- 2013 Less Rhyme, More Reason: Knowledge-based Poetry Generation with Feeling, Insight and Wit. In *Proc. of ICCO 2013, the 4th Int. Conference on Computational Creativity*. Sydney, Australia.

Tony Veale

- 2014 Running With Scissors: Cut-Ups, Boundary Friction and Creative Reuse. In *Proceedings of ICCBR-2014, the 22nd International Conference on Case-Based Reasoning*.

Eileen Cornell Way

- 1991 *Knowledge Representation and Metaphor: Studies in Cognitive systems*. Kluwer.

Joseph Weizenbaum

- 1966 ELIZA – A Computer Program For the Study of Natural Language Communication Between Man And Machine. *Communications of the ACM* 9 (1): 36–45.

Yorick Wilks

- 1978 Making Preferences More Active. *Artificial Intelligence* 11(3):197-223.