

Argument Box

Rehaf AlJammaz,¹ Yasheng She,² Michael Mateas³

¹raljamma@ucsc.edu

²yashe@ucsc.edu

³michaelm@soe.ucsc.edu

University of California, Santa Cruz

Abstract

Games currently feature simple models of morality and moral reasoning. These typically take the form of a simple binary opposition between good and evil, in which game actions are sorted into these categories, and reputation systems that track the reputation of players (and sometimes NPCs) with in-game factions based on the actions taken. This paper presents a more sophisticated model of moral reasoning based on Lakoff’s metaphorical family models: the strict father and nurturant parent morality. This model utilizes more rapid, surface level categorization of different situations in moral categories, and deeper reasoning that characterizes situations in terms of their relationship to moral virtues as determined by the metaphor value system. This model is used in the experimental prototype Argument Box (AB), a social argument simulator where the player argues with clients visiting the shop. Arguments in AB center around the moral virtues and vices of simulated characters in the social simulation Talk of the Town. This paper presents the current architecture, discussing its technical details.

Introduction

In many faction and reputation-based games, morality is often portrayed as a simple binary value of good and evil. The reputation of the player character (and sometimes the Non-Player Character (NPC)) is represented as a single-dimensional scale where players or NPCs land on a spectrum based on their deeds and the morality points they accumulate.

Unfortunately, NPCs in these types of games tend to employ surface-level judgments about the characters they encounter. For instance, games such as World of Warcraft and Neverwinter Nights (Blizzard 2004; ObsidianEntertainment 2002) employ NPCs that attack opposing factions or views without any reasoning from the characters involved, usually based on high-level decisions such as the character belonging to an opposing faction. These simplistic moral decisions and binary systems can affect a character’s believability, making them feel mechanical, stereotyped, and less life-like. By developing comprehensive moral models, we can

create characters with better behavior understandability, that are less predictable, and have more sophisticated awareness of different social situations. All of which are dimensions of believability (Gomes et al. 2013). How a character reasons about these moral situations can enhance the illusion of believability; one aspect of reasoning involves belief modeling.

When referring to belief modeling in games, we usually associate it with facts and theories of mind. For example, a given character might believe that an NPC has brown eyes or another NPC likes them. Research games employ this factual sense of belief modeling, such as characters believing in false information in (Ryan et al. 2015) and (Guimaraes, Santos, and Jhala 2017). Beliefs under this category have also been used alongside social systems (Evans and Short 2013; Morais, Dias, and Santos 2019), creating unique stories. They have even been used as game mechanics, such as the player influencing NPC actions by directly ‘injecting’ false beliefs in MKULTRA (Horswill 2015). In this paper, we refer to a broader definition of belief modeling that includes an NPC’s moral beliefs and values.

We believe there is a need to create a more sophisticated moral reasoning system. This paper presents our current progress in Argument Box: a game prototype that employs moral reasoning for NPCs in terms of both surface and deep values.

Our prototype adds a layer of moral reasoning by incorporating Lakoff’s metaphorical family models: the strict father and nurturant parent morality into our agent’s belief systems. In the current prototype, we focus on the Strict Father Model (SFM), which, as the name suggests, (Lakoff 2010) employs a metaphorical strict father figure as the head of the household, guiding our agent’s actions. The SFM (Lakoff 2010) believes that “the world is a dangerous place.” Therefore the “parent” needs to protect the children, become authoritative and teach them “right” from “wrong,” reward them if they do good, and punish them if they do wrong. Each family system employs a number of metaphors or values they believe in; for example, the SFM values strength. Anything that reduces strength is seen as weak and immoral.

Our model classifies different situations according to each character’s held moral beliefs at a surface level. Surface-level reasoning is used to initially describe an NPC’s stance

on an issue. We also incorporated deeper reasoning capabilities that are referenced when an NPC strongly cares for a given topic, granting an NPC the ability to fight for a passionately held belief when the player disagrees with them. Our deeper model allows for deeper reasoning that defines NPC situations related to moral values as determined by the father model metaphor.

Related Work

In this section, we review related work in terms of morality, reasoning, and social simulation as it relates to various elements of our project.

Social Simulation and Belief-Based Reasoning

Several research games incorporate social simulation and reasoning in their system’s designs. Work such as PromWeek (McCoy et al. 2012) and successor CiF-based systems (Morais, Dias, and Santos 2019; Guimaraes, Santos, and Jhala 2017) utilize social physics as the main mechanic of their design, influencing player conversations and NPC interactions through social elements such as an NPC’s relationship with other NPCs.

Unsurprisingly, social simulation and belief modeling are closely related to one another. CiF-CK (Guimaraes, Santos, and Jhala 2017) adapted the CiF architecture (McCoy et al. 2010) by incorporating belief modeling alongside their social networks, making it possible for characters to believe in false information. Comme il Faut - Exiles (CiF-EX) (Morais, Dias, and Santos 2019) extended belief modeling beyond the interacting character, where other NPCs can infer relationships between characters in the world.

In Talk of the Town (TotT) (Ryan et al. 2015), a historical town simulator, incorporates characters that can forget information, misremember and lie. Beliefs in TotT are affected by a character’s social network, what characters remember, and the strength of the information supporting a given belief. One interesting application of TotT is Bad News (Samuel et al. 2016), where players converse with real actors portraying different NPCs. The actor acts on instructions from a developer behind the scenes. We incorporate TotT’s generated characters as topics of conversation in Argument Box.

Other work, such as MKULTRA (Horswill 2015), focused on belief modeling, where beliefs are ‘injected’ into an NPC’s knowledge base, allowing the player to manipulate the NPC directly. In Versu’s (Evans and Short 2013) model of belief, agents can share public views of the world and have specific instances of individual false beliefs creating unique gameplay experiences.

While social networks do not directly influence this project, we use characters that incorporate social relationships and elements in their design. We essentially used Talk of the Town’s (Ryan et al. 2015) exported characters as conversation topics between the player and the NPC.

Moral Reasoning

Morality in Games There are many games that involve morality and moral decision-making. In commercial games,

moral decision-making is usually tied to the player’s actions. Games such as The Witcher 3, Undertale, and those by Telltale Games (CDProjekt 2015; Fox 2015; TellTaleGames 2004) place the player in high stakes-moral situations where their action affects other NPCs and the story’s outcome; this is usually limited to the player character, without modeling NPC rationalization in moral dilemmas. Other games that do model NPC moral decision-making do so at somewhat superficial levels. For instance, NPCs in World of Warcraft (Blizzard 2004) and Farcry (LCGEntertainment 2004) use high-level faction judgments, such as attacking NPCs in opposing factions on sight.

In academia, morality systems in games often target moral theories or thought experiments. Works such as Togelius (Togelius 2011) and Nelson (Nelson 2012) examine Kant’s categorical imperative, which states, “Act only according to that maxim by which you can at the same time will that it should become a universal law.” Togelius’s prototype involves a procedural system that defines a new maximum rule each time an event triggers, following the categorical imperative principle. Nelson’s prototype started with game rules and allowed the player to break them, with the game adding new rules making the rule breaking universal. According to Nelson (Nelson 2012), this quickly led to havoc. Togelius arrived at a similar conclusion noting the difficulty of designing meaningful gameplay around universalization of player action.

Other authors have built systems based on political divides as well as belief modeling; Azad and Martens’s system, Lyra (Azad and Martens 2019; 2018) simulates character interactions and conversations in politically charged groups. What is unique about Azad et al.’s system is their character’s ability to learn new biases and introduce it to their knowledge bases.

Lakoff’s Moral Politics Lakoff is well known for his work on metaphor theory, where he identifies deep metaphors, such as *Life is a Journey* or *Love is War*, that structure human cognition and the use of language (Lakoff and Johnson 2008). In more recent work, Lakoff has proposed two governing metaphors underlying political and moral reasoning. These two metaphors, the *Strict Father* and *Nurturant Parent*, govern the conservative and liberal (respectively) world views regarding the individual and their relationship to society. Our work on moral reasoning for NPCs described in this paper builds on this work.

Lakoff’s (Lakoff 2010) strict father model employs a metaphorical strict father as the head of a figurative household. The strict father is seen as an authoritative figure that values strength, discipline and believes in rewards and punishments. The strict father views morality in terms of one’s ability to abide by certain values such as strength, moral order (following the natural order of things), and Moral Boundaries (i.e., deviating from the norm is wrong). Anything that reduces these values is seen as immoral. Lakoff illustrates an example where the act of purchasing illegal drugs is seen as immoral, as it emerges from low self-control (i.e. low strength) according to SFM.

On the other hand, the Nurturant Parent believes that

members can grow as a result of nurturance and care. It encourages self-reliance by caring for others. Unlike the SFM, the figurative parents in the NPM value their member’s opinions; the model employs respect between its members rather than dreading punishments or expecting rewards. Like the SFM, the NPM has metaphors and values it believes in, such as Morality as Fair Distribution, Morality as Social Nurturance (e.g., strengthening social relationships with others and mending those relationships), and Morality as Nurturance (i.e., being regularly empathetic).

In our current prototype, we focus on the values employed by the Strict Father Model (SFM). The deep values held by this model include Moral Strength, Moral Boundaries, and Moral Wholeness. We will further describe how we incorporated these metaphors into our system in a later section.

The prototype and Gameplay

In a classic Monty Python skit called ‘Argument Clinic’ (Monty Python), a man approaches a clinician asking to buy an argument; he then proceeds to argue with the clinician about arguments, debating if their argument is an argument! Likewise, our prototype features an argument simulator game, where NPCs walk into a shop called Argument Box (AB) to procure arguments with the player character, focusing on arguments about moral behavior.

The prototype described in this paper starts when an NPC initiates conversations about characters living in their town. They give the player a piece of gossip about a particular character they heard about, including their opinion and evaluation of that character’s actions. The player can agree, disagree or vouch for the talked-about character by arguing for particular stances. Depending on the player’s response and how passionately the NPC feels about the current argument, the NPC can reference their deep-seated beliefs as moral arguments, which for the presented prototype are based on Lakoff’s strict father model. In the following sections, we will discuss the system in greater detail. First, we examine the components of our system at a high level; we then look at our NPC’s architecture, covering its creation, design, and belief modeling. Lastly, we go through two examples covering a typical argument in our simulation. We note that the characters in Argument Box are modeled as polyhedra. Thus, in the example dialog below, there are references to “shapes” which in our world is the same thing as saying “people.”

High Level Overview

Our current prototype features a character coming into a clinic-like shop called Arg Box to argue with the player character about the latest gossip in town. These characters are called conversational NPCs (CNPC). The conversational NPCs (CNPC) are the main NPCs the player converses within the Argument Box.

A single game loop involves the CNPC conversing about one of the town’s characters, stating what they heard, and saying how they feel about it. The conversed-about characters are what we refer to in this paper as the background NPCs (BNPC). The BNPCs include a separate list of NPCs that the player character never interacts with; they are used

to seed conversational topics by our CNPCs. These BNPCs are generated from TotT. TotT, as a reminder, is a historical town simulator that generates characters, including various elements such as character relationships, locations, and jobs.

We import the data from TotT in JSON format. We then search for patterns, such as combinations of attributes or social connections on BNPCs or temporal sequences undergone by BNPCs, in a process similar to story sifting (Kreminski, Dickinson, and Wardrip-Fruin 2019; Max Kreminski 2021). BNPCs are assigned tags based on the patterns that matched. We then filter the BNPCs by thresholding the number of tags (only BNPCs with enough tags are potential topics of conversation) and place the filtered list of BNPCs into a priority queue based on the number of tags found and the quality of the tags. Tags that involve multi-character patterns are weighted more highly than tags resulting purely from within-character patterns. We provide examples of the sifting patterns we use in a later section.

Once we have our list of BNPCs, the CNPC brings up appropriate dialogue based on the tags present on the BNPCs. The starting dialogue is unbiased and simply expands the tag into a textual utterance. For example, the tag *familyPerson* gets translated into “Have you heard that X has a large family?”

Each pattern in the system is further mapped to the CNPC’s surface values. These surface values denote how a CNPC generally feels about a set of tags; there is a many-to-many mapping of tags to surface values which is described in more detail below. The CNPC holds each surface value with a strength of high, medium or low, indicating how passionately the CNPC believes in that value.

As long as the player agrees with the judgements of the CNPC, the conversation will stay at the level of surface values. However, when the player disagrees with the CNPC on a surface value they hold strongly, the CNPC will perform deeper reasoning using the strict father model to back up their claim. The conversational options provided to the player allow them to agree, disagree, bring up a specific discussion, or change the conversation topic entirely if they desire. Figure 1 illustrates the general components of the system. The following sections will describe how each component works in greater detail.

Modeling BNPCs

BNPCs are generated using the TotT simulation. Our sifting patterns focus on topics we can have moral debates about, and so generally exclude details such as locations of homes and businesses, street names and physical character descriptions. We currently use 58 patterns to assign tags to BNPCs as follows:

- Patterns that directly map a single TotT raw attribute. Examples include *isWealthy*, *departed*, and *familyPerson*, indicating if a character has wealth, left town, or has a family, respectively.
- Patterns that are created by combining different TotT raw attributes. For example, a retired character that is 45 years old has the pattern *retiredYoung* assigned.

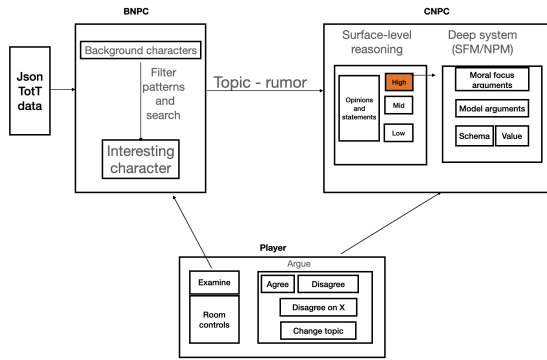


Figure 1: Overview of the system, featuring the components of the player, BNPC and CNPC architectures

- Patterns that are based on TotT character jobs; these patterns are used to make controversial assumptions about characters, based on the character’s career and its potential effect on others. For example, any TotT character with the job *miner* or *cooper* is assigned the tag *polluterRole*. This is later used for conversations related to the environment.
- Patterns that focus on relationships with other generated characters, such as love triangles and backstabbing. Examples include *friendWithBestFriendsEnemy* and *InLoveWithSpouseOfFriend*.
- Patterns that combine lower-level tags into higher-level ones. For example, if a BNPC has the tags *adultButNot-working* and *IsWealthy* (from lower-level patterns) this results in the higher-level tag *notWorkingAndRich*.
- The last category of patterns combines tags found by other patterns with the Five Favor Personality traits (Goldberg 1990; Ryan et al. 2015) that TotT assigns to characters. For example, the pattern *hasA lotOfEnemies* combined with the personality factor of high Agreeableness results in the tag *tooTrustingOfEnemies*.

As mentioned above, once all pattern matching is completed and the tags have been assigned, BNPCs are placed in a priority queue based on the number and quality of tags. By starting the conversation about characters at the front of the queue, this ensures that the CNPC will have a good number of debatable topics to argue about with the player.

Modeling CNPCs

Surface Values

The CNPC starts the conversation by picking the BNPC at the front of the queue to talk about, doing this until it has exhausted all the BNPCs tags or the player chooses another character to talk about.

Once a tag has been chosen, the CNPC starts the conversation by commenting on the tag in a value neutral manner. For example, if our BNPC named Mike had the tag *familyPerson* selected by the system, the CNPC states, “Oh, have you heard that Mike has a big family?” The CNPC then states

how they feel about this tag by relating it to their surface values.

There are currently 28 surface values defined in our system, of which a subset will be held (with varying strengths of low, medium or high) by a CNPC. Examples include *LoveIsForFools*, *LoveAboveAllElse*, *FamilyPerson*, and *ShapeIsNothingIfNotSocial*. BNPC tags map in a many-to-many way with surface values. The tags that map to a given surface value are called core tags of that surface value. The core tags have a many-to-many relationship with the surface values. For example, the tag *willActOnLove* is a core tag of both the surface values *BeTrueToYourHeart* and *LoveIsForFools*. In the case that a core tag maps to two or more mutually held surface values, this provides some non-determinism on how the CNPC will comment on the presence of this tag, depending on which surface value is taken as being activated.

When a CNPC is instantiated, the system randomly assigns the surface values and their accompanying strength. Some surface values are mutually exclusive, so can not be simultaneously held with high strength. For example, if our CNPC holds *LoveIsForFools* with a high rating, it cannot hold *BeTrueToYourHeart* with any strength other than low.

Additionally, the mapping from core tags to surface values is used for the conversational options presented to the player. This allows the player to bring up BNPC characteristics during the conversation that explicitly agree or disagree with the CNPC at the level of surface values.

The surface values are used for immediate value-laden reactions during the conversation. As long as the player and CNPC agree with each other, the conversation can stay at the surface value level. However, when the player’s judgments disagree with the CNPC, the system switches to reasoning about deep values as determined by the Strict Father Model (SFM). Implementing a version of the Nurturant Parent Model (NPM) is ongoing work. Switching to this deeper model allows the CNPC to marshal arguments by bringing up characteristics that relate to more deeply held values. This prevents the conversation from immediately degenerating into repeated assertions (e.g. “Yes it is! Not it isn’t! Yes it is! Not it isn’t”) at the surface level.

Deep Values: Strict Father Model

We define six deep values drawn from Lakoff’s book *How Liberals and Conservatives think* (Lakoff 2010) to specify the SFM:

- *Moral Boundaries* warns about the danger of deviating from the norm. Characters that deviate from the norm are seen as immoral by a character holding the SFM.
- *Self Interest* sees seeking one’s self-interest as moral and interfering with one’s self-interest as immoral.
- *Moral Wholeness* is concerned with unity and conformity among characters.
- *Moral Essence* evaluates a character’s past actions as indicators for their future actions, making the assumption that past actions are the result of a character’s “essence”.
- *Strength* values a character’s ability to act in or handle difficult or sensitive situations. Low strength is seen as

Sample: surface to deep argument structure -
disagreement High - SFM

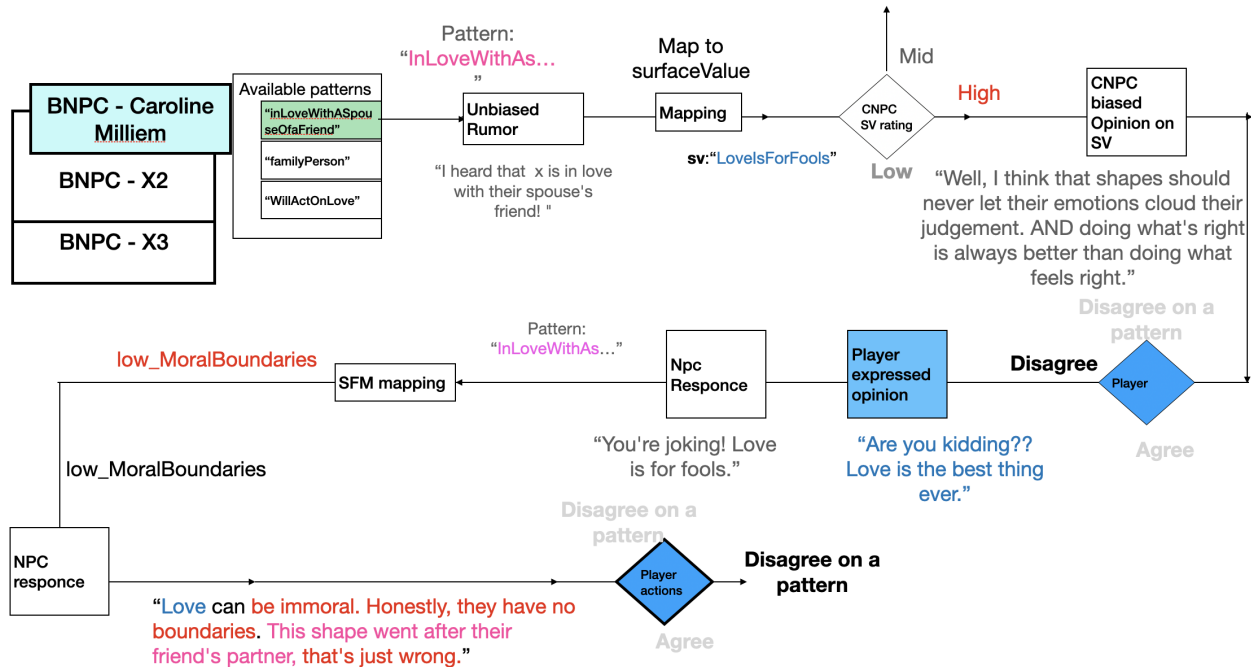


Figure 3: Deep value conversation loop - player disagrees with high surface value

BNPC have been explored; if they haven't been explored, the system proceeds with the same loop but with the newly appointed tag *familyPerson*. Otherwise, the system moves on to the next BNPC in the queue.

A Deep rooted conversation - Using the SFM

Let us take the previous example but assume our CNPC actually cares deeply about the surface value *LovesForFools*. The CNPC conveys this as, "Well, I think that shapes should never let their emotions cloud their judgment AND doing what's right is always better than doing what feels right." As a note on authoring, we use AND and BUT to emphasize the CNPC's high or low stances by extending the same sentence with a modifier that enhances how strongly or weakly they feel about a given topic. This helps minimize the combinatorial amount of dialog we have to write and signals the underlying model more strongly to the player.

In this example the player disagrees with the CNPC, selecting the dialog option: "Are you kidding?? Love is the best thing ever." Again, the tone mimics that of the CNPC. The CNPC then responds with, "You're joking! Love is for fools." It consults the underlying model, in this case the SFM, to back up why love is indeed for fools.

The currently selected tag *inLoveWithSpouseOfFriend* is validated as a core tag of *LovesForFools*, and determined to score low against the deep value *Moral Boundaries*. Thus the SFM validates the surface value stance, and allows the CNPC to make an argument based on deep values: "Love can be immoral. Honestly, they have no boundaries. This shape went after their friend's partner; that's just wrong."

We note that the text is written in a way that references the specific tag, illustrates the deep value, and provides reasoning as to why the surface value *LovesForFools* should be held. The player can then agree or disagree with the pattern, choose another BNPC pattern to bring up, passing it as an argument for the current surface value.

If the CNPC is presented with a tag that the SFM maps in a way contradictory to the surface value argument being made, it then searches through the BNPC's available tags for an alternative argument. If any of the remaining tags are core tags of the surface value and the SFM mapping supports the argument, it presents it as a backup argument. Otherwise, the CNPC mimics a person backed into a corner, randomly firing off defenses based on the tags found. This results in more generic arguments, for example, stating "But that BNPC has a family!" in response to the situation where the BNPC is in love with the spouse of a friend.

Current Limitations and Future Work

This paper presents our current work on Argument Box, a game prototype incorporating a moral reasoning system that operates on both surface and deep values (beliefs). The model is inspired by Lakoff's work on family-based metaphors in moral reasoning. As future work we are implementing the Nurturant Parent Model, refining the conversation loop so that under appropriate (but challenging) circumstances the player can change the CNPCs mind, and playtesting the game.

References

- Azad, S., and Martens, C. 2018. Addressing the elephant in the room: Opinionated virtual characters. In *AIIDE Workshops*.
- Azad, S., and Martens, C. 2019. Lyra: Simulating believable opinionated virtual characters. In *Proceedings of the AAAI conference on artificial intelligence and interactive digital entertainment*, volume 15, 108–115.
- Blizzard. 2004. World of warcraft. CD-ROM.
- CDProjekt. 2015. The witcher 3: Wild hunt. Xbox One.
- Evans, R., and Short, E. 2013. Versu—a simulationist storytelling system. *IEEE Transactions on Computational Intelligence and AI in Games* 6(2):113–130.
- Fox, T. 2015. Undertale. PC.
- Goldberg, L. R. 1990. An alternative” description of personality”: the big-five factor structure. *Journal of personality and social psychology* 59(6):1216.
- Gomes, P.; Paiva, A.; Martinho, C.; and Jhala, A. 2013. Metrics for character believability in interactive narrative. In *International conference on interactive digital storytelling*, 223–228. Springer.
- Guimaraes, M.; Santos, P.; and Jhala, A. 2017. Cif-ck: An architecture for social npcs in commercial games. In *2017 IEEE Conference on Computational Intelligence and Games (CIG)*, 126–133. IEEE.
- Horswill, I. D. 2015. Mkultra. In *Eleventh Artificial Intelligence and Interactive Digital Entertainment Conference*.
- Kreminski, M.; Dickinson, M.; and Wardrip-Fruin, N. 2019. Felt: a simple story sifter. In *International Conference on Interactive Digital Storytelling*, 267–281. Springer.
- Lakoff, G., and Johnson, M. 2008. *Metaphors we live by*. University of Chicago press.
- Lakoff, G. 2010. *Moral politics: How liberals and conservatives think*. University of Chicago Press.
- LCGEntertainment. 2004. Far cry primal. PC.
- Max Kreminski, Melanie Dickinson, M. M. 2021. Winnow: A domain-specific language for incremental story sifting. In *Conference on Artificial Intelligence and Interactive Digital Entertainment*.
- McCoy, J.; Treanor, M.; Samuel, B.; Tarse, B.; Mateas, M.; and Wardrip-Fruin, N. 2010. Comme il faut 2: A fully realized model for socially-oriented gameplay. In *Proceedings of the Intelligent Narrative Technologies III Workshop*, 1–8.
- McCoy, J.; Treanor, M.; Samuel, B.; Reed, A. A.; Wardrip-Fruin, N.; and Mateas, M. 2012. Prom week. In *Proceedings of the International Conference on the Foundations of Digital Games*, 235–237.
- Monty Python, u. c. Argument - monty python.
- Morais, L.; Dias, J.; and Santos, P. A. 2019. From caveman to gentleman: a cif-based social interaction model applied to conan exiles. In *Proceedings of the 14th International Conference on the Foundations of Digital Games*, 1–11.
- Nelson, M. 2012. Prototyping kant-inspired reflexive game mechanics. In *Proceedings of the 2012 Workshop on Research Prototyping in Games*.
- ObsidianEntertainment, B. 2002. Neverwinter nights. [CD-ROM].
- Ryan, J.; Summerville, A.; Mateas, M.; and Wardrip-Fruin, N. 2015. Toward characters who observe, tell, misremember, and lie. In *Proceedings of the AAAI Conference on Artificial Intelligence and Interactive Digital Entertainment*, volume 11.
- Samuel, B.; Ryan, J.; Summerville, A. J.; Mateas, M.; and Wardrip-Fruin, N. 2016. Bad news: An experiment in computationally assisted performance. In *International Conference on Interactive Digital Storytelling*, 108–120. Springer.
- TellTaleGames. 2004. Telltalegames. PC.
- Togelius, J. 2011. A procedural critique of deontological reasoning. In *DiGRA Conference*.