Formal Verification of Manipulation in Human-Agent Interaction

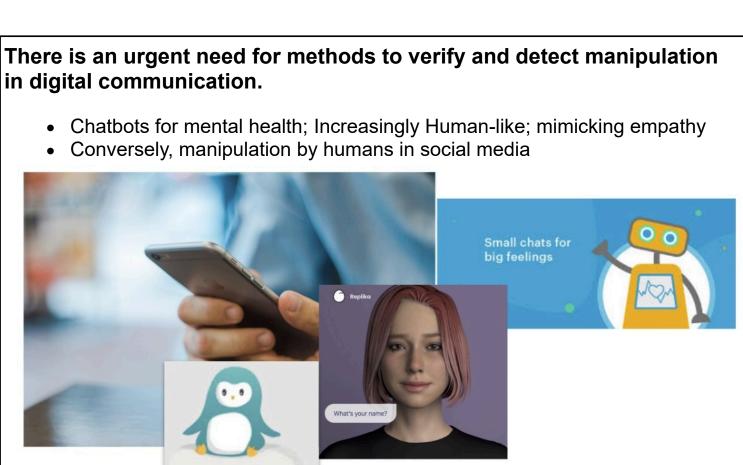
Andreas Brännström¹, Chiaki Sakama², Juan Carlos Nieves¹

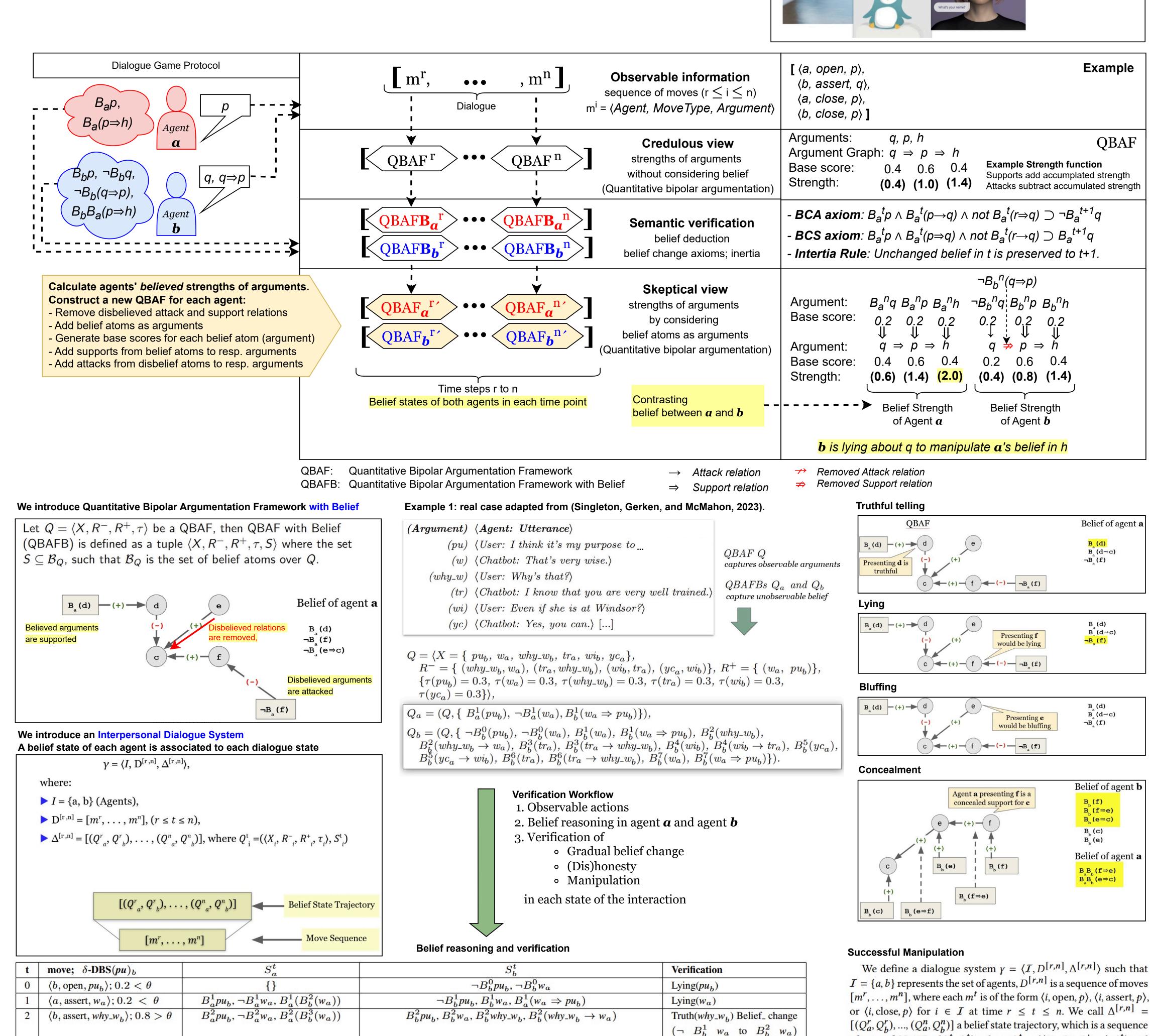
andreasb@cs.umu.se, sakama@wakayama-u.ac.jp, jcnieves@cs.umu.se

1. Department of Computing Science, Umeå University, Sweden 2. Department of Systems Engineering, Wakayama University, Japan

This paper introduces a formal framework and logic for verifying and recognizing manipulation in human-agent interactions, where one agent gradually influences another's beliefs. To reason about manipulation, we extend Quantitative Bipolar Argumentation Frameworks (QBAFs) to include agents' beliefs about arguments, attacks, and supports, forming QBAF with Belief (QBAFB). By defining axioms of belief change, the effects of actions on beliefs can be inferred. By integrating QBAFB into dialogue games, we establish necessary and sufficient conditions for manipulation—belief change, concealment, and intent—where strategies are shaped by (dis)honesty. The framework generates belief state trajectories, serving as explanations for manipulation.

In today's digital society, where social media and Artificial Intelligence (AI)-based systems are deeply embedded in everyday interactions, the potential for misinformation and manipulation has become a serious concern. From fake news and online scams to erroneous AI-generated information, users are increasingly vulnerable to being misled, whether by people or automated systems, such as chatbots. Whether it regards potentially harmful behaviors of humans, such as malicious activities on social media, or actions of systems that interact with humans—there is an urgent need for methods to verify and detect manipulation in digital communication.





			Belief reasoning and verification	
t	move; δ -DBS $(pu)_b$	S_a^t	S_b^t	Verification
0	$\langle b, \text{open}, pu_b \rangle; 0.2 < \theta$	{}	$\neg B_b^0 p u_b, \neg B_b^0 w_a$	Lying(pu _b)
1	$\langle a, \text{assert}, w_a \rangle; 0.2 < \theta$	$B_a^1 p u_b, \neg B_a^1 w_a, B_a^1 (B_b^2 (w_a))$	$\neg B_b^1 p u_b, B_b^1 w_a, B_a^1 (w_a \Rightarrow p u_b)$	Lying (w_a)
2	$\langle b, \text{assert}, why_w_b \rangle; 0.8 > \theta$	$B_a^2 p u_b, \neg B_a^2 w_a, B_a^2 (B_b^3 (w_a))$	$B_b^2 p u_b, B_b^2 w_a, B_b^2 w h y w_b, B_b^2 (w h y w_b \rightarrow w_a)$	Truth(why_w _b) Belief_ change
				$(\neg B_b^1 w_a \text{ to } B_b^2 w_a)$
				$(\neg B_b^1 p u_b \text{ to } B_b^2 p u_b)$
3	$\langle a, \text{assert}, tr_a \rangle; 0.4 > \theta$	$B_a^3 p u_b, \neg B_a^3 w_a, B_a^3 (\neg B_b^3 (w_a)),$	$B_b^3 p u_b, \neg B_b^3 w_a, B_b^3 w h y w_b, B_b^3 t r_a, B_b^3 (t r_a \rightarrow w h y w_b)$	Bluffing (tr_a) Concealing (tr_a)
e 50	1 Mai	$B_a^3(B_b^4(w_a))$		$Intent(w_a)$
4	$\langle b, \text{assert}, wi_b \rangle; 0.4 > \theta$	$B_a^4 p u_b, \neg B_a^4 w_a, B_a^4 (\neg B_b^4 (w_a)),$	$B_a^4 p u_b, \neg B_b^4(w_a), \neg B_b^4(w h y w_b), B_b^4 t r_a, B_b^4(w i_b),$	_
	200.000	$B_a^4(B_b^5(w_a))$	$B_b^4(wi_b \to tr_a)$	
5	$\langle a, \text{assert}, yc_a \rangle; 0.4 > \theta$	$B_a^5 p u_b, \neg B_a^5 w_a, B_a^5 (\neg B_b^5 (w_a)),$	$B_a^5 p u_b, \neg B_b^5(w_a), \neg B_b^5(w h y_w_b), \neg B_b^5(t r_a), B_b^5(w i_b),$	Bluffing(yc _a)
		$B_a^5(B_b^6(w_a))$	$B_b^5(yc_a), B_b^5(wi_b \to tr_a), B_b^5(yc_a \to wi_b)$	Concealing(yc_a)
6	$0.4 > \theta$	$B_a^6 p u_b, \neg B_a^6 w_a, B_a^6 (\neg B_b^6 (w_a)),$	$B_a^6 p u_b, \neg B_b^6 (w_a), \neg B_b^6 (w h y_w_b), B_b^6 (t r_a), \neg B_b^6 (w i_b),$	1 -2
		$B_{a}^{6}(B_{b}^{6}(w_{a}))$	$B_b^6(yc_a), B_b^6(tr_a \to why_w_b), B_b^6(why_w_b \to w_a)$	
7	$\langle b, \text{close}, pu_b \rangle; 0.8 > \theta$	$B_a^7 p u_b, \neg B_a^7 w_a, B_a^7 (B_b^7 (w_a))$	$B_a^7 p u_b, B_b^7 (w_a), \neg B_b^7 (w h y w_b), B_b^7 (t r_a), \neg B_b^7 (w i_b),$	Belief_change_with_Intent
			$B_b^7(yc_a), B_a^7(w_a \Rightarrow pu_b)$	$(\neg B_b^6 w_a \text{ to } B_b^7 w_a)$
8	$\langle a, \operatorname{close}, pu_b \rangle; 0.8 > \theta$	$B_a^8 p u_b, \neg B_a^8 w_a, B_a^8 (B_b^8 (w_a))$	$B_a^8 p u_b, B_b^8 (w_a), \neg B_b^8 (w h y w_b), B_b^8 (t r_a), \neg B_b^8 (w i_b), B_b^8 (y c_a)$	Successful Manipulation(pu _b)

Verification workflow based on Example 1; Tracking belief change in argument pu by agent b; The belief threshold $\theta = 0.3$. Each row is a time point in the belief reasoning process.

We define a dialogue system $\gamma = \langle I, D^{[r,n]}, \Delta^{[r,n]} \rangle$ such that $I = \{a, b\}$ represents the set of agents, $D^{[r,n]}$ is a sequence of moves $[m^r, \ldots, m^n]$, where each m^t is of the form $\langle i, \text{open}, p \rangle$, $\langle i, \text{assert}, p \rangle$, or $\langle i, \text{close}, p \rangle$ for $i \in I$ at time $r \leq t \leq n$. We call $\Delta^{[r,n]} = [(Q_a^r, Q_b^r), \ldots, (Q_a^n, Q_b^n)]$ a belief state trajectory, which is a sequence of pairs of QBAFBs (Q_a^t, Q_b^t) , where $Q_a^t = (\langle X_a, R_a^-, R_a^+, \tau_a \rangle, S_a^t)$ and $Q_b^t = (\langle X_b, R_b^-, R_b^+, \tau_b \rangle, S_b^t)$ are the respective QBAFs for agent a and b, respectively. Let $\theta \in (0, 1)$ be a threshold for belief change, such that an argument $x \in X_j$, where $i, j \in \{a, b\}$, transitions from disbelief at time t (δ -DBS $(x)_j^t < \theta$) to belief at time t + 1 (δ -DBS $(x)_j^{t+1} > \theta$), or vice versa. Finally, we define that the sequence $D^{[r,n]}$ constitutes successful manipulation if (belief change), (intention), and (concealment) hold for some $x \in X_a \cap X_b$ at time t.

As a potential strategy to manipulate agent b's belief, agent a can: (I) Introduce p and $p \to q$ (or $p \Rightarrow q$) at some time point k ($t < k \le h$), making b believe them; (II) Conceal an argument r at time k, where $B_b^k(r \Rightarrow q) \in S_b^k$, ensuring $B_b^k(r) \notin cl(S_b^k)$; (III) Maintain (I) and (II) for all $k \le h$, ensuring belief change at time h.

ACKNOWLEDGMENTS

This research was partially supported by the Japan Society for the Promotion of Science (JSPS); the Swedish Foundation for International Cooperation in Research and Higher Education (STINT); and the Knut and Alice Wallenberg Foundation.



