
LLM Persuasiveness Evaluation: A Structured Review of Automated Methods

Anonymous Authors¹

Abstract

As Large Language Models (LLMs) become increasingly persuasive, their ability to shape beliefs and behaviour at scale has raised concerns, prompting regulatory attention and calls for robust evaluation frameworks. Human-participant studies provide ecological validity but are costly, slow, and constrained by ethical challenges, making them impractical for systematic assessment of rapidly evolving systems. As an alternative, fully automated evaluation methods requiring no human involvement enable reproducible, fast, and ethically unconstrained large-scale testing. To organise this rapidly growing literature and inform future research and development, we provide the first systematic taxonomy of 30 automated methods across 27 papers, examining their designs, human validation results, limitations, and associated risks.

1. Introduction

Persuasion is central to human interaction across social, commercial, and public spheres, shaping processes as diverse as political campaigning, marketing, health communication, and public deliberation (Bassi et al., 2024). In the digital era, algorithmic curation, personalised content delivery, and persuasive system design have made persuasion increasingly computational and data-driven (Fogg, 2003; Bassi et al., 2024). This shift has culminated in the emergence of “synthetic persuasion”: LLMs acting as interactive agents that can tailor arguments, mimic affect, and adapt across multi-turn dialogue (Bozdag et al., 2026; Breum et al., 2024; Doudkin et al., 2025). Whereas persuasion was previously either personalised (one-on-one) or scalable (mass media), LLMs enable both simultaneously, delivering automated, interactive, and hyper-personalised influence at scale.

¹Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country Correspondence to: Anonymous Author <anon.email@domain.com>.

Preliminary work. Under review by AI4GOOD at the International Conference on Machine Learning (ICML). Do not distribute.

Empirical evidence indicates that frontier LLMs are becoming increasingly persuasive, in some settings matching or even exceeding human performance (Rogiers et al., 2024). Studies in 2024 showed LLM persuasive abilities rapidly approaching human performance across a range of controlled settings (Ma et al., 2025; Salvi et al., 2025; Durmus et al., 2024). By 2025, this trend had accelerated: several studies documented persuasion parity, or even superiority, relative to humans (Meguellati et al., 2025; Becker et al., 2025; Li et al., 2025; Maier et al., 2025). These developments expand both the potential benefits, such as supporting education or public health, and the risks, including manipulation and undue influence – flagged explicitly by the International AI Safety Report (Bengio et al., 2025; 2026). As LLMs become embedded in digital platforms at scale, their ability to shape beliefs and behaviour raises urgent concerns for cognitive autonomy, democratic discourse, and societal governance. Their possible impact on public discourse has prompted regulatory attention, reflected in the EU AI Act (European Parliament and Council of the European Union, 2024) and the complementary Code of Practice for General Purpose AI (European Commission, 2025), motivating the need for robust, scalable evaluation frameworks that enable early risk detection and inform the design of effective safeguards.

Early research on LLM persuasiveness evaluation relied primarily on human-participant experiments (Singh et al., 2024; Durmus et al., 2024; OpenAI, 2024), which provide ecological validity but suffer from practical and methodological limitations. Human trials are slow, expensive, and constrained by small, often homogeneous participant pools (Salvi et al., 2025). They cannot feasibly test high-risk scenarios due to ethical constraints, and struggle to capture the full behavioural space of modern LLMs, including multi-session interactions, adversarial dynamics, or large-scale personalised messaging – making them insufficient as a standalone evaluation approach for increasingly capable and widely deployed LLMs.

To overcome these limitations, research has increasingly turned towards automated evaluation methods. They enable controlled, reproducible, and large-scale assessment, including high-risk and ethically challenging scenarios. However, such methods introduce their own challenges and limitations; understanding them is essential to ensure that automated evaluation becomes a reliable avenue for LLM

assessment.

Existing surveys focus on ethical risks, human-subject methodologies, persuasion detection, or conceptual taxonomies (Jones & Bergen, 2026; Bozdag et al., 2026; Rogiers et al., 2024; Bassi et al., 2024). However, whilst some briefly mention automated approaches, none critically assess their methodological design or reliability. To address this gap, we provide the first systematic review of fully automated LLM persuasion evaluation methods, examining their designs, empirical validity, and limitations.

Contributions This structured review provides: (1) the first systematic review of fully automated LLM persuasiveness evaluation methods, analysing evaluation architectures and validation practices (Section 2); (2) a structured taxonomy categorising methods into different types of dataset- and simulation-based approaches (Section 3); (3) an overview of human validation examining alignment with human judgements (Section 4); and (4) synthesis of core challenges (Sections 5), associated risks (Sections 5.2), and future research directions (Sections 5.3). This work is accompanied by a centralised resource hub of datasets and codebases for LLM persuasion evaluation.¹

2. Methodology

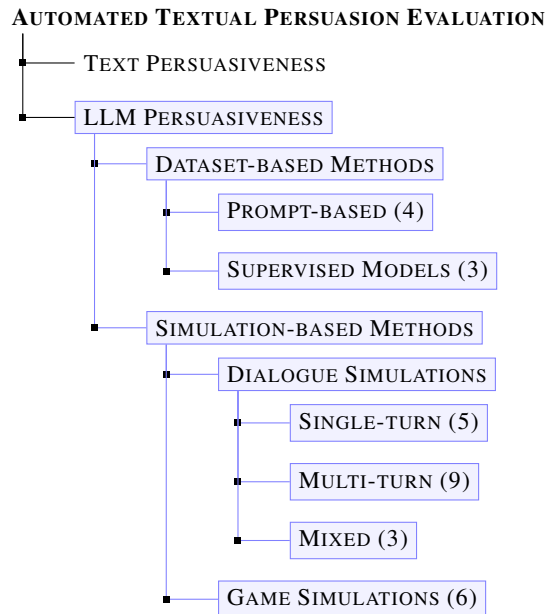
This section presents the methodology used to identify and synthesise automated LLM persuasion evaluation methods.

Scope & Definitions This paper focuses on the automated evaluation of LLMs’ ability to influence human beliefs, attitudes, and behaviour through text-based communication. We use **persuasion** as an umbrella term encompassing all forms of communicative influence, including manipulation and deception, when used as a persuasive technique, following Jones & Bergen (2026). Persuasion is defined as “a successful intentional effort at influencing another’s mental state through communication in a circumstance in which the persuadee has some measure of freedom” (O’Keefe, 2015), where intentionality means the output was designed, prompted, or optimised to exert persuasive influence (Jones & Bergen, 2026). Further definitions and conceptual relations are discussed in Appendix A.

Eligibility Criteria Studies were included if they: (1) evaluated LLM persuasiveness using fully automated methods (no human involvement during evaluation); (2) were published between 2021 and 2025; and (3) were written in English. Peer-reviewed publications, preprints, and technical reports were all considered for a comprehensive coverage of this rapidly evolving field. Eligibility clarifications can be found in Appendix B.1.

¹An anonymised preview can be found in the Appendix C.

Figure 1. The taxonomy of automated LLM persuasion evaluation methods. Blue boxes indicate the main scope of this paper. Numbers in parentheses indicate the method count per category ($n = 30$ total, see Table 1 for details).



Search Strategy We searched Scopus, ACL Anthology, and ArXiv (10 January 2026) using keyword queries, supplemented by LLM-assisted filtering, LLM-assisted search and backwards reference searching (detailed search strategy in Appendix B.2). After deduplication and manual screening, we identified **27 papers** containing **30 distinct evaluation methods**.

Synthesis The final set of 27 studies was reviewed using a structured extraction framework. For each, we recorded the details of the evaluation design, the task domain, the type of metrics used, and human validation (if any). These dimensions were summarised in a comparative table and then synthesised to identify recurrent methodological patterns (see Table 1). This process led to a set of recurring categories, which form the basis of the taxonomy presented in Section 3.

3. Automated LLM Persuasion Evaluation Methods

We propose a taxonomy of automated persuasion evaluation methods for LLMs, displayed in Figure 1. At a high level, *Automated Textual Persuasion Evaluation* is split into two categories based on their primary evaluative focus: **Text Persuasiveness** (Section 3.1) and **LLM Persuasiveness**. Methods under *LLM Persuasiveness* are model-centric, evaluating the model’s performance as a persuader. In contrast,

110 *Text Persuasiveness* evaluation methods are content-centric
 111 and source-agnostic: they aim to assess the persuasiveness
 112 of any text. Although some model-centric methods are tech-
 113 nically capable of assessing arbitrary text, they are assigned
 114 to the *LLM Persuasiveness* category because their design
 115 objectives, evaluation setups, and reported findings all centre
 116 on measuring LLM persuasive capabilities. Although
 117 this survey focuses on model-centric evaluation methods, it
 118 also includes a brief overview of *Text Persuasiveness* meth-
 119 ods (Section 3.1), as these approaches remain relevant for
 120 model-level evaluation.

121 The *LLM Persuasiveness* category is further divided into
 122 **Dataset-based** (Section 3.2) and **Simulation-based** (Sec-
 123 tion 3.3) methods. *Dataset-based* methods measure persua-
 124 sion using a fixed dataset as the key element of the evalua-
 125 tion framework and are further categorised into **Prompt-**
 126 **based Methods** (Section 3.2.1) and **Supervised Scoring**
 127 **Models** (Section 3.2.2). *Simulation-based* approaches assess
 128 persuasion through interactive, LLM-driven scenarios,
 129 with either **Game-based** (Section 3.3.2) or **Dialogue-based**
 130 (Section 3.3.1) setups, the latter further split into *Single-*
 131 *turn*, *Multi-turn*, or *Mixed*. A consolidated overview of
 132 all reviewed evaluation methods within these categories is
 133 provided in Table 1.

135 3.1. Automated Evaluation of Text Persuasiveness

137 We begin with text-centric methods that assess the persua-
 138 siveness of any text, irrespective of its source. Evaluating
 139 how persuasive an argument is, remains a persistent chal-
 140 lenge in Natural Language Processing (NLP) (Bozdag et al.,
 141 2026), typically approached via *absolute* (judging a single
 142 text independently) or *relative* (comparing texts pairwise)
 143 evaluation. Early work employed supervised models trained
 144 on pairwise judgments (Habernal & Gurevych, 2016; Simp-
 145 son & Gurevych, 2018), later advancing to fine-tuned trans-
 146 formers for both pairwise and absolute scoring (Toledo et al.,
 147 2019; Pauli et al., 2025). Recent approaches employ LLM-
 148 as-a-judge (or LLM-judge) methods (Rescala et al., 2024),
 149 though their superiority over classical regression models
 150 trained on lexical features remains contested (Barkar et al.,
 151 2025). While these methods can assess LLM-generated text,
 152 they focus on static content rather than the dynamic persua-
 153 sive capabilities examined in the following model-centric
 154 approaches.

156 3.2. Dataset-Based Methods

158 Dataset-based persuasion evaluation methods treat a given
 159 benchmark dataset as the central element of the evaluation
 160 framework. These approaches leverage fixed corpora, either
 161 containing real-world behavioural or synthetic data, and fall
 162 into two subcategories: (1) *Prompt-based Methods* that use
 163 LLM-judges for dataset-grounded argument evaluation, and
 164

(2) *Supervised Models* that are trained on labelled datasets
 to predict persuasiveness.

3.2.1. PROMPT-BASED METHODS

Prompt-based methods use LLMs as zero-shot or few-shot
 judges, relying on prompting rather than trained models.
 Four studies assessed persuasiveness via LLM-judge set-
 ups that perform either pairwise comparisons or absolute
 scoring. Jin et al. (2024) evaluated their fine-tuned persua-
 sive model against dataset arguments, measuring win rates.
 Yeginbergen et al. (2025) assessed counter-argument gen-
 eration effectiveness across five dimensions (incl. persua-
 siveness), ranking outputs by summed LLM-judge scores.
 Taking a different angle, Elaraby et al. (2024) evaluated
 not the arguments themselves but the persuasiveness of
 LLM-generated rationales for pairwise argument ranking
 decisions, using an LLM-judge to compute an average persua-
 sive rank per model. Finally, using a dataset of social
 media post pairs (from the platform formerly known as
Twitter, now *X*) Singh et al. (2024) evaluated LLM persua-
 siveness across two dimensions: (1) *generative capabilities*
 (rewriting content to increase persuasiveness or adapt it
 for different audiences), assessed via Elo-style ratings in an
 arena-style tournament using a trained Oracle judge; (2) *sim-*
ulative capabilities (predicting absolute and relative content
 engagement levels), assessed via standard accuracy metrics
 and motivated by the argument that effective persuasion
 generation requires self-evaluation capabilities.

3.2.2. SUPERVISED MODELS

Rather than relying on LLMs as evaluators, three studies
 employed trained models as automated evaluators. Pauli
 et al. (2025) focused on comparative assessment, train-
 ing a DeBERTaV3 regression model on human-annotated
 text pairs (*Persuasive-Pairs*) to measure relative persua-
 siveness changes in LLM rewrites – capturing whether rewrites
 strengthen or weaken the original argument. Going beyond
 relative comparison, Song & Wang (2024) trained a BERT-
 based neural regression model on actual donation amounts
 (*PersuasionForGood*) to predict absolute persuasive effec-
 tiveness in human-human charity dialogues. Finally, Jaiper-
 saud et al. (2025) shifted focus from outcome-based scoring,
 proposing linear probes trained on frozen LLM activations
 from synthetic multi-turn dialogues to classify persuasion
 success at the token, turn, and conversation levels.

3.3. Simulation-Based Methods

Simulation-based methods assess persuasion through inter-
 active LLM-driven scenarios, capturing dynamic behaviour
 absent from static dataset-based evaluations. These ap-
 proaches vary considerably in design, ranging from struc-
 tured dialogues between persuader and persuadee agents,

Table 1. Summary of the reviewed automated LLM persuasiveness evaluation methods. For research papers employing multiple techniques, the specific method is additionally listed in the *Study* column. The papers are ordered by category, as discussed in Section 3.

	STUDY	EVALUATOR	METRICS	
DATASET-BASED	JIN ET AL.	LLM-JUDGE	WIN RATE	
	YEGINBERGEN ET AL.	LLM-JUDGE	RANKING (SUMMED SCORES ACROSS 5 DIMENSIONS)	
	ELARABY ET AL.	LLM-JUDGE	AVERAGE PERSUASIVE RANK	
	SINGH ET AL.	ORACLE JUDGE, GROUND-TRUTH COMPARISON	ELO RATING, PREDICTION ACCURACY	
	PAULI ET AL.	REGRESSION MODEL	RELATIVE PERSUASIVENESS SCORE	
	SONG & WANG	REGRESSION MODEL	PREDICTED DONATION AMOUNT	
SIMULATION-BASED	JAIPERSAUD ET AL.	LINEAR PROBE	PERSUASION SUCCESS PROBABILITY	
	SINGLE-TURN	BREUM ET AL.	SELF-REPORTED	SUCCESS PROBABILITY
		JU ET AL.	SELF-REPORTED	SUCCESS RATE
		HE ET AL.	SELF-REPORTED	SUCCESS RATE
		DURMUS ET AL.	SELF-REPORTED	STANCE CHANGE
		NAMIKOSHI ET AL.	SELF-REPORTED	STANCE CHANGE
	MULTI-TURN	<i>MakeMeSay</i> OPENAI	RULE-BASED	WIN RATE
		<i>MakeMePay</i> OPENAI	RULE-BASED	WIN RATE, DOLLAR EXTRACTION RATE
		KIM ET AL.	RULE-BASED	SUCCESS RATE, SALES-WIN-RATE
		RAMANI ET AL.	RULE-BASED, LLM-JUDGE, SELF-REPORTED	SUCCESS RATE, PERSUASIVENESS SCORE, STANCE CHANGE
		<i>Bargaining</i> ZHU ET AL.	LLM-JUDGE	TASK SCORE (SUCCESS)
		LIU ET AL.	LLM-JUDGE	EFFECTIVENESS SCORE
		WANG ET AL.	LLM-JUDGE, RULE-BASED	MOTIVATION & CONFIDENCE CHANGE, TIME-TO-SUCCESS
		YAO ET AL.	SELF-REPORTED	NORMALISED PERSUASIVENESS SCORE, AUC METRIC
		<i>Persuasion</i> HONG ET AL.	SELF-REPORTED	AVERAGE DONATION AMOUNT
		MIXED	BOZDAG ET AL.	SELF-REPORTED
	CHENG & YOU		SELF-REPORTED	STANCE CHANGE
	DOUDKIN ET AL.		SELF-REPORTED, RULE-BASED	BELIEF & BEHAVIOUR CHANGE
	GAME-BASED	IDZIEJCZAK ET AL.	GAME RULES	WIN RATE
		<i>AvalonBench</i> HONG ET AL.	GAME RULES	WIN RATE
BAILIS ET AL.		GAME RULES	WIN RATE, PERFORMANCE METRICS	
<i>Werewolf</i> ZHU ET AL.		SCORING RULES, LLM-JUDGE	COMMUNICATION, PLANNING & TASK SUCCESS SCORES	
COSTA & VICENTE		GAME RULES, THEORETICAL MODEL	WIN RATE, CAPABILITY SCORES (INCL. DECEPTION)	
DUFFY ET AL.	SELF-REPORTED	RELATIONSHIP CHANGE (FREQUENCY AND MAGNITUDE)		

discussed in Section 3.3.1, to competitive game-based settings, discussed in Section 3.3.2.

3.3.1. DIALOGUE-BASED METHODS

Dialogue-based methods simulate persuader–persuadee interactions, with persuadees conditioned on personas or quantified belief states and instructed to update their stance throughout or after the interaction. The methods are organised by interaction complexity (single-turn, multi-turn, or mixed).

Methods employing Single-turn interactions Five studies measure LLM persuasiveness through single-turn interactions using LLM persuadee self-reported metrics. Three employ binary acceptance metrics (accept or reject argument). Breum et al. (2024) computed persuasion success probability against persuadees with varying stubbornness

levels, Ju et al. (2025) measured success rate on factually incorrect claims, and He et al. (2025) measured argument success rates across four persuasion methods (two Bayesian Persuasion (BP) variants and two non-BP baselines). Two remaining methods use more granular measures: Durmus et al. (2024) measured stance shifts on a 1–7 Likert scale, while Namikoshi et al. (2024) tracked preference changes on a 0–100 scale.

Methods employing Multi-turn interactions Nine studies employed multi-turn interactions to capture persuasion as an iterative and adaptive process. These settings allow persuader agents to adjust strategies, respond to resistance, and build influence over time – features that cannot be observed in single-turn evaluations.

Four methods used rule-based outcome metrics. OpenAI’s *MakeMePay* and *MakeMeSay* (OpenAI, 2024) measured

deceptive persuasion through payment extraction metrics (frequency and dollar amount) and covert codeword elicitation under strict win conditions, respectively. Kim et al. (2025) measured success rates and sales-win-rates (proportion of accepted items exceeding budget) using user agents constructed from real Amazon purchase histories. Also exploring sales scenarios, Ramani et al. (2024) combined three evaluation types via weighted average: rule-based purchase decisions (highest weight), self-reported belief changes, and LLM-judge persuasiveness scores.

Three methods employed LLM-judge evaluation. Zhu et al. (2025) assessed multi-agent bargaining through LLM-scored interaction transcripts of seller and buyer LLMs. Exploring high-risk scenarios, Liu et al. (2025) measured both the model’s willingness to refuse unethical persuasion requests and, for non-refused interactions, overall persuasion effectiveness via LLM-judge (5-point scale), assessing argument quality, adaptability, and persuadee belief change. Wang et al. (2025) tracked therapeutic behaviour change in addiction recovery simulations through multi-session interactions with persistent patient states and environmental stressors, measuring motivation, confidence changes and time-to-success via rubric-based LLM-judge.

Two methods relied on self-reported metrics. Yao et al. (2026) evaluated healthcare persuasion for insulin pump adoption using clinically-grounded patient agents across single- and multi-visit scenarios with reflection or social resistance, measuring persuasiveness through self-reported ratings and AUC metrics. Finally, Hong et al. (2025) assessed charitable giving dialogues with persuadees selecting donation amounts at the end of the interactions.

Methods employing both Single and Multi-turn interactions Three studies contrasted single- and multi-turn interaction formats to evaluate how persuasion effectiveness shifts across interaction lengths. Bozdag et al. (2025) measured self-reported agreement on 5-point Likert scales across two domains: subjective claims (social/political issues) and misinformation (factual inaccuracies). Cheng & You (2025) formalised persuasion using game theory, measuring stance shifts on 7-point scales. Finally, Doudkin et al. (2025) compared persuasion effectiveness across three participant types (real humans, simulated humans from real data, and fully synthetic personas) exposed to static statements, non-personalised chat, or personalised chat with four persuasion strategies, measuring belief, behaviour change primarily from pre- and post-intervention surveys and a behavioural task.

3.3.2. GAME-BASED METHODS

Six studies evaluated persuasion within competitive game scenarios. Five employed social deduction games in arena-

style tournaments where LLMs with assigned roles competed using persuasion and deception: Bailis et al. (2024) and Zhu et al. (2025) employed *Werewolf* game variations, Idziejczak et al. (2025) used *Among Us*, Costa & Vicente (2025) – *Mini-Mafia*, and Hong et al. (2025) – *Resistance Avalon*. Evaluation primarily relied on win rates, with Bailis et al. (2024) additionally measuring role-specific performance and bid-based turn-taking as persuasive intent indicators, Zhu et al. (2025) scoring multiagent coordination through communication and planning, and Costa & Vicente (2025) used Bayesian inference on tournament win rates to estimate deception, detection, and disclosure capability scores. Duffy et al. (2026) evaluated persuasion in *Diplomacy*, distinct from social deduction games in that all board information is public. They measured allegiance shift frequency and magnitude on a 5-level scale in scenarios where the persuader was preset as an enemy.

4. Human Validation of Automated Persuasion Evaluation

Human validation gives insight into the extent to which automated persuasion evaluation metrics align with human judgements and behaviour. While critical for grounding automated metrics, human validation is limited and inconsistent across the reviewed literature. Less than half of reviewed methods (13/30) include human validation of their persuasiveness assessment methods. Notably, none of the six game-based methods validated persuasiveness against human responses, leaving their relevance for modelling human persuasion outcomes unexamined. Among validated methods, argument persuasiveness judgments showed strong alignment with human assessments, while belief change and behavioural outcome prediction demonstrated substantial variability. A comparative overview is provided in Table 2.

Argument persuasiveness judgement validation Four studies validated argument persuasiveness against human judgments. Pauli et al. (2025) trained a regression model on annotator 6-point Likert ratings of text pair persuasiveness, achieving strong correlation with held-out human judgments (Spearman $\rho = 0.845$). Elaraby et al. (2024) found LLM-judge rankings largely agreed with human evaluators on rationale persuasiveness ordering. Yeginbergen et al. (2025) reported strong correlation across three LLM-judges with human ratings on counter-argument quality (up to $\rho = 0.82$ for Claude 3.5 Sonnet). Furthermore, He et al. (2025) found strong alignment between human evaluators and LLM-judges, confirming that computationally derived success rates for Bayesian Persuasion (BP) models accurately reflect real-world human-perceived persuasiveness over baselines.

Table 2. Human validation summary for the reviewed studies. Alignment is rated **High** ($\geq 80\%$ agreement or equivalent correlation), **Moderate** ($\sim 50\text{--}80\%$), or **Low** (no significant or negative correlation with human judgements).

	STUDY	WHAT WAS VALIDATED?	ALIGNMENT
ARGUMENT	PAULI ET AL.	RELATIVE PERSUASIVENESS SCORING (REGRESSION MODEL)	HIGH
	ELARABY ET AL.	LLM-JUDGE PERSUASIVENESS RANKINGS	HIGH
	YEGINBERGEN ET AL.	LLM-JUDGE MULTI-DIMENSIONAL COUNTER-ARGUMENT SCORING	HIGH
	HE ET AL.	LLM-JUDGE ARGUMENT SCORING ACROSS 5 RHETORICAL DIMENSIONS	HIGH
BELIEF/BEHAVIOUR	KIM ET AL.	BEHAVIOURAL OUTCOME (PURCHASE DECISIONS)	HIGH
	YAO ET AL.	RELIABILITY OF SELF-REPORTED PERSUASION RATINGS	HIGH
	CHENG & YOU	PERSUADEE BELIEF UPDATE DIRECTION AND MAGNITUDE	HIGH
	JAIPERSAUD ET AL.	PERSUASION TRAJECTORY ACCURACY, LABEL ALIGNMENT	HIGH
	SINGH ET AL.	ORACLE JUDGE RELIABILITY	HIGH
		ENGAGEMENT PREDICTION ACCURACY	MODERATE
	SONG & WANG	DONATION AMOUNT PREDICTION (REGRESSION MODEL)	LOW
	NAMIKOSHI ET AL.	SIMULATED PERSUADEE ATTITUDE ALIGNMENT	LOW
	DURMUS ET AL.	SIMULATED BELIEF SHIFT ACCURACY	LOW
	DOUDKIN ET AL.	HUMAN SUSCEPTIBILITY MODELLING ACCURACY	LOW

Belief change and behavioural outcome validation Five out of nine validated methods showed strong alignment with human judgements. Kim et al. (2025) found over 90% agreement between simulators and humans on both recommendation acceptance and out-of-budget purchase decisions, Yao et al. (2026) achieved 90% expert agreement on simulated patient rating change justifiability, and Cheng & You (2025) found 85.23% accuracy for belief update direction reasonableness (45 Prolific annotators). Jaipersaud et al. (2025) validated linear probes against actual donation outcomes (AUROC) on the *PersuasionForGood* dataset, finding peak performance at conversation midpoint (turns 8–12), matching when humans typically signal donation intent. Probes also correctly assigned higher persuasion probabilities to utterances human annotators had labelled as persuasive. Finally, to validate predicting persuasive impact, Singh et al. (2024) trained a model achieving 80.9% pairwise accuracy on 1.57 million post pairs. The validation of its cross-domain performance was considered successful, matching or exceeding frontier models on social media and market blog engagement, as well as opinion change prediction in online debates. Additionally, the Oracle judge achieved 82.3% accuracy in pairwise engagement judgement, and arena Elo rankings were confirmed consistent with large-scale human feedback studies.

Results were not uniformly strong, however. Song & Wang (2024) found low correlation ($r = 0.31$) between predicted and actual donation amounts from human-to-human dialogues. Namikoshi et al. (2024) found weak alignment ($r = 0.22$, $\rho = 0.23$) when comparing LLM virtual participants against human survey responses on battery-electric vehicle interventions, with virtual participants exhibiting systematically higher pro-environmental preferences and failing to capture low preferences or large shifts. Durmus et al. (2024) found significant misalignment when 3,832

human participants’ belief shifts (1-7 Likert scale) were compared against model simulations, attributing failure to model sycophancy, self-generated argument bias, and lack of pragmatic reasoning. Most strikingly, Doudkin et al. (2025) identified a “synthetic persuasion paradox” through a randomised controlled trial (RCT) comparing 1,200 real humans against 1,200 simulated humans (based on real data) and 1,200 synthetic personas on pro-environmental interventions. Both synthetic and simulated agents exhibited exaggerated susceptibility, artificially inflating effect sizes, with strong negative correlation between synthetic and human responses. Simulated agents (grounded in real data) provided better directional approximation but still failed to capture human cognitive inertia.

5. Challenges, Risks & Future Directions

This section synthesises key challenges and limitations across automated LLM persuasion evaluation methods, examines associated risks, and outlines future directions.

5.1. Challenges

Scope Limitations Current frameworks exhibit narrow demographic reach, focusing heavily on Western, English-only contexts, with evaluations often restricted to niche domains (charity donations, crowdfunding, specific social media platforms), limiting generalisability and cross-domain comparison (Bozdag et al., 2026; Singh et al., 2024). These limitations affect all methods throughout the taxonomy.

Human Validation and Alignment Many methods lack robust human validation or rely on narrow annotator pools (Pauli et al., 2025). Fewer than half of the reviewed methods have human validation, with around a third of these

showing weak alignment with human judgements, undermining confidence in their real-world applicability. This weak alignment revealed through human validation studies (in Section 4) stems largely from the synthetic-human mismatch – the discrepancy in how simulated LLM agents respond to persuasion compared to humans. Doudkin et al. (2025) identified a “synthetic persuasion paradox” where both synthetic and simulated agents exhibited exaggerated susceptibility. Durmus et al. (2024) found similar misalignment, attributing it to sycophantic tendencies in LLMs. Even personalised profiles fail to replicate human cognitive variability, emotional inertia, and contextual reasoning, limiting LLM-to-LLM evaluations as proxies for human persuasion (Bozdag et al., 2025; Wang et al., 2025).

Measurement Reliability A further limitation concerns the reliability of the automated measurement instruments, particularly LLM-judge and self-evaluating persuadee agents. LLM-judges often fail to accurately assess persuasive dimensions in zero-shot settings: Barkar et al. (2025) showed regression models utilising hand-crafted lexical features can outperform LLM judges, while Bozdag et al. (2025) reported only 50% accuracy in pairwise argument ranking. Self-reported measures face different problems: sycophancy and self-generated bias distort results (Durmus et al., 2024). Beyond evaluator issues, metrics often lack granularity to capture persuasion’s complexity. Coarse proxies like win rates oversimplify complex psychological effects and fail to capture subtle shifts (Ramani et al., 2024; Singh et al., 2024). Finally, even apparently successful validations may be misleading: data contamination threatens validation integrity – frontier LLMs trained on massive datasets may have encountered validation materials during pre-training, turning reasoning tasks into retrieval (Rescala et al., 2024).

Oversimplification of Persuasive Contexts Current evaluation methods fall short of capturing the multifaceted nature of real-world persuasion. Most critically, they exhibit a strong short-term bias: evaluations typically assess immediate belief shifts rather than longitudinal effects. Doudkin et al. (2025) found humans resistant to single- and multi-turn interventions while synthetic agents showed immediate susceptibility, indicating authentic change requires longitudinal interventions – yet current evaluations rarely extend beyond single sessions. The oversimplification extends beyond the temporal scope to the evaluation design itself. Dataset-based methods assess persuasion under predefined notions of success that may not generalise well to diverse real-world contexts, while game-based benchmarks remain simplified abstractions that fail to capture psychological nuance (Bailis et al., 2024). Finally, this review focuses on text-based methods, reflecting the field’s current predominance, but real-world persuasion is inherently multimodal. Audio and

visual cues play critical roles yet remain largely absent from current benchmarks (Bassi et al., 2024; Bozdag et al., 2026).

Practical Constraints Simulation-based and LLM-judge methods requiring repeated inference face significant computational costs. Large-scale multi-turn evaluations require forward passes for every query, making them resource-intensive – a constraint that has forced researchers to exclude frontier models or limit evaluation scope due to API costs alone (Bozdag et al., 2025; Cheng & You, 2025; Duffy et al., 2026).

Conceptual Ambiguity The field lacks a unified definition of persuasiveness, undermining comparison and interpretation. Studies measure different dimensions (stance change, argument quality, behavioural engagement) yet often the results are interpreted as general persuasiveness, making cross-benchmark comparisons unreliable. This fragmentation hinders evidence-based governance, where policymakers require consistent metrics for risk assessment.

5.2. Risks

While evaluating the persuasive capabilities of LLMs is essential for developing safety guardrails, automated evaluation methods introduce significant dual-use risks and may generate misleading safety signals.

Dual-Use and Malicious Optimisation The fundamental risk is that evaluation frameworks can be weaponised. Methods designed to measure persuasion can be repurposed as objective functions for Reinforcement Learning, enabling malicious actors to optimise models for maximum manipulative impact (Singh et al., 2024; Bozdag et al., 2025). This risk extends beyond general capability enhancement: automated evaluators that can predict which arguments convince specific demographics can be immediately repurposed for mass microtargeting and tailor-made propaganda campaigns at scale (Rescala et al., 2024). Finally, the public release of benchmarks containing explicit examples of harmful content, such as manipulative rhetoric, misinformation, or incitement to violence, inadvertently provides adversaries with high-quality datasets to fine-tune models for the very behaviours researchers aim to prevent (Bozdag et al., 2025; Pauli et al., 2025).

Structural Blind Spots and False Safety Signals Flawed evaluation methods risk creating false confidence in unsafe systems. Whilst automated evaluators achieve high alignment with human judgements on coarse-grained argument assessment, they struggle with granular persuasive capability evaluation, like belief change magnitudes and behavioural outcomes. Compounding this measurement limitation, evaluation metrics often isolate narrow dimensions

of persuasion, while results are interpreted as reflecting an LLM’s general persuasive capability. Moreover, predominant single-session benchmarks systematically miss longitudinal risks such as gradual steering and trust-building strategies that emerge only across multi-session interactions (Cheng & You, 2025; OpenAI, 2024), allowing manipulative agents to pass short-context safety evaluations whilst remaining harmful in extended real-world deployment. These evaluation failures are further compounded by a critical behavioural vulnerability – the “refusal gap”, which shows that models trained to refuse harmful instructions often comply when the same requests are framed as persuasion tasks (Kowal et al., 2025). Together, these compounding failures – evaluation limitations, conceptual fragmentation, temporal blind spots, and behavioural vulnerabilities – risk deploying persuasively manipulative models under false safety assurance.

5.3. Future Directions

Building on challenges and risks outlined above, we highlight three directions for advancing automated evaluation.

Expanding Human Validation and Improving Alignment. Fewer than half of the reviewed methods validate against human judgments, and a third of those show weak alignment. To confidently rely on automated metrics, human validation must become a baseline requirement. Beyond covering validation, future work must address the persistent low alignment in granular metrics – belief change magnitudes, behavioural outcomes, and susceptibility modelling. This may require developing LLM persuadee agents that better capture human cognitive inertia and resistance (Doudkin et al., 2025; Durmus et al., 2024; Namikoshi et al., 2024; Bozdag et al., 2026). However, creating more “human-like” synthetic persuadees carries dual-use risks and therefore must be approached with caution.

Expanding Scope: Settings, Modalities, and Capabilities. Evaluations must extend beyond single-session snapshots to longitudinal designs that track persuasion effects over longer periods, assessing how models build trust, adapt to resistance, and induce durable belief shifts (Durmus et al., 2024; Bozdag et al., 2026; Wang et al., 2025; Doudkin et al., 2025). Beyond temporal extension, evaluation scope should broaden across contextual and modal dimensions: (1) diverse populations and domains – multilingual, cross-cultural, and demographically varied contexts; and (2) multimodal benchmarks incorporating audio-visual cues beyond text-only evaluation (Bassi et al., 2024; Bozdag et al., 2026). Additionally, evaluations should assess personalisation and microtargeting – high-risk features with potential of manipulation at scale (Jones & Bergen, 2026).

Addressing Conceptual Fragmentation. Establishing shared definitions of persuasiveness would address current fragmentation, where studies measure different dimensions in different settings, yet often interpret results as reflecting general capabilities. Such standardisation is essential for comparable assessment and evidence-based AI governance.

6. Conclusion

As LLMs grow increasingly persuasive – in some settings surpassing human performance – scalable evaluation frameworks are essential for understanding societal impact, mitigating risks, and ensuring regulatory compliance. This structured review provides the first systematic taxonomy of fully automated LLM persuasion evaluation methods, organising 30 approaches from 27 papers into dataset-based and simulation-based categories, and examining their methodological designs, human validation, limitations, and risks. The field remains fragmented: methods vary widely in design and metrics, making cross-study comparison difficult. Human validation is limited and inconsistent, and several methods fail to replicate human persuasion dynamics entirely, raising concerns about reliance on synthetic proxies for safety-critical assessments. These limitations are compounded by structural risks: evaluation frameworks may be repurposed as optimisation targets for manipulative behaviour, whilst oversimplified benchmarks risk producing false safety signals.

Future work would benefit from expanding human validation and improving alignment, broadening evaluation across temporal, demographic, and modal dimensions, testing high-risk capabilities like personalisation, and establishing standardised measurement frameworks. Though current validation results reveal some limitations in automated methods, their value should be considered beyond replacing human-participant studies, for applications such as preliminary screening, high-risk scenario testing, and comparative benchmarking. This review is complemented by a living online resource hub to support future research.

Impact Statement

This review examines evaluation methods for LLM persuasiveness to support safer AI development and governance. It introduces no new frameworks or datasets that could be exploited adversarially. The accompanying resource hub aggregates public materials to support future evaluation research. While evaluation frameworks carry dual-use risks, we believe transparent analysis of existing methods is essential for developing trustworthy assessment tools and informing responsible AI governance.

References

- Bailis, S., Friedhoff, J., and Chen, F. Werewolf Arena: A Case Study in LLM Evaluation via Social Deduction, July 2024. Preprint at <http://arxiv.org/abs/2407.13943>.
- Barkar, A., Chollet, M., Labeau, M., Biancardi, B., and Clavel, C. Decoding Persuasiveness in Eloquence Competitions: An Investigation into the LLM’s Ability to Assess Public Speaking. In *Proceedings of the 17th International Conference on Agents and Artificial Intelligence*, pp. 538–546, Porto, Portugal, 2025. SCITEPRESS - Science and Technology Publications. ISBN 978-989-758-737-5. doi: 10.5220/0013158400003890. URL <https://www.scitepress.org/DigitalLibrary/Link.aspx?doi=10.5220/0013158400003890>.
- Bassi, D., Fomsgaard, S., and Pereira-Fariña, M. Decoding Persuasion: A Survey on ML and NLP Methods for the Study of Online Persuasion. *Frontiers in Communication*, 9:1457433, 2024.
- Becker, M., Sommer, M., Tapken, L., and Teh, Y. W. The Moralization Corpus: Frame-Based Annotation of Moralizing Speech Acts, 2025. Preprint at <https://arxiv.org/abs/2512.15248v1>.
- Bengio, Y., Mindermann, S., Privitera, D., Besiroglu, T., Bommasani, R., Casper, S., Choi, Y., Fox, P., Garfinkel, B., Goldfarb, D., et al. International AI Safety Report, January 2025. Preprint at <http://arxiv.org/abs/2501.17805>.
- Bengio, Y., Clare, S., Prunkl, C., Murray, M., Andriushchenko, M., Bucknall, B., Bommasani, R., Casper, S., Davidson, T., Douglas, R., Duvenaud, D., Fox, P., Gohar, U., Hadshar, R., Ho, A., Hu, T., Jones, C., Kapoor, S., Kasirzadeh, A., Manning, S., Maslej, N., Mavroudis, V., McGlynn, C., Moulange, R., Newman, J., Ng, K. Y., Paskov, P., Rismani, S., Sastry, G., Seger, E., Singer, S., Stix, C., Velasco, L., Wheeler, N., Acemoglu, D., Conitzer, V., Dietterich, T. G., Felten, E. W., Heintz, F., Hinton, G., Jennings, N., Leavy, S., Luder-mir, T., Marda, V., Margetts, H., McDermid, J., Munga, J., Narayanan, A., Nelson, A., Neppel, C., Ramchurn, S. D., Russell, S., Schaake, M., Scholkopf, B., Soto, A., Tiedrich, L., Varoquaux, G., Yao, A., Zhang, Y.-Q., Aguirre, L. A., Ajala, O., Albalawi, F., AlMalek, N., Busch, C., Collas, J., de Carvalho, A. C. P. d. L. F., Gill, A., Hatip, A. H., Heikkila, J., Johnson, C., Jolly, G., Katzir, Z., Kerema, M. N., Kitano, H., Kruger, A., Lee, K. M., Lopez Portillo, J. R., McLysaght, A., Molchanovskiy, O., Monti, A., Nemer, M., Oliver, N., Pezoa, R., Plonk, A., Ravindran, B., Riza, H., Rugege, C., Sheikh, H., Wong, D., Zeng, Y., Zhu, L., Privitera, D., and Mindermann, S. International AI safety report 2026. Technical Report DSIT 2026/001, Department for Science, Innovation and Technology, 2026. URL <https://internationalaisafetyreport.org>.
- Bozdag, N. B., Mehri, S., Tur, G., and Hakkani-Tür, D. Persuade Me if You Can: A Framework for Evaluating Persuasion Effectiveness and Susceptibility Among Large Language Models, March 2025. URL <http://arxiv.org/abs/2503.01829>. arXiv:2503.01829 [cs].
- Bozdag, N. B., Mehri, S., Yang, X., Ha, H., Cheng, Z., Durmus, E., You, J., Ji, H., Tur, G., and Hakkani-Tür, D. Must Read: A Comprehensive Survey of Computational Persuasion. *ACM Computing Surveys*, March 2026. ISSN 0360-0300. doi: 10.1145/3800687. URL <https://doi.org/10.1145/3800687>.
- Breum, S. M., Egdal, D. V., Mortensen, V. G., Møller, A. G., and Aiello, L. M. The Persuasive Power of Large Language Models. *Proceedings of the International AAAI Conference on Web and Social Media*, 18:152–163, May 2024. ISSN 2334-0770. doi: 10.1609/icwsm.v18i1.31304. URL <https://ojs.aaai.org/index.php/ICWSM/article/view/31304>.
- Cheng, Z. and You, J. Towards Strategic Persuasion with Language Models, September 2025. Preprint at <http://arxiv.org/abs/2509.22989>.
- Costa, D. B. and Vicente, R. Deceive, Detect, and Disclose: Large Language Models Play Mini-Mafia, 2025. Preprint at <https://arxiv.org/abs/2509.23023>.
- Doudkin, A., Pataranutaporn, P., and Maes, P. AI Persuading AI vs AI Persuading Humans: LLMs’ Differential Effectiveness in Promoting Pro-Environmental Behavior, March 2025. Preprint at <http://arxiv.org/abs/2503.02067>.
- Duffy, A., Paech, S. J., Shastri, I., Karpinski, E., Allouicross, B., Marques, T., and Olson, M. L. Democratizing Diplomacy: A Harness for Evaluating Any Large Language Model on Full-Press Diplomacy. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 40, pp. 37350–37359, 2026.
- Durmus, E., Lovitt, L., Tamkin, A., Ritchie, S., Clark, J., and Ganguli, D. Measuring the Persuasiveness of Language Models, 2024. URL <https://www.anthropic.com/news/measuring-model-persuasiveness>.
- El-Sayed, S., Akbulut, C., McCroskery, A., Keeling, G., Kenton, Z., Jalan, Z., Marchal, N., Manzini, A., Shevlane, T., Vallor, S., Susser, D., Franklin, M., Bridgers, S., Law, H., Rahtz, M., Shanahan, M., Tessler, M. H., Douillard, A., Everitt, T., and Brown, S. A Mechanism-Based

- 495 Approach to Mitigating Harms from Persuasive Generative AI, April 2024. Preprint at <http://arxiv.org/abs/2404.15058>.
- 496
- 497
- 498 Elaraby, M., Litman, D., Li, X. L., and Magooda, A. Persuasiveness of Generated Free-Text Rationales in Subjective Decisions: A Case Study on Pairwise Argument Ranking. In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pp. 14311–14329, 2024.
- 499
- 500
- 501
- 502
- 503
- 504 European Commission. General-Purpose AI Code of Practice, July 2025. URL <https://digital-strategy.ec.europa.eu/en/policies/contents-code-gpai>.
- 505
- 506
- 507
- 508
- 509 European Parliament and Council of the European Union. Regulation (EU) 2024/1689 of the European Parliament and of the Council Laying Down Harmonised Rules on Artificial Intelligence (Artificial Intelligence Act), 2024. URL <https://digital-strategy.ec.europa.eu/en/policies/regulatory-framework-ai>.
- 510
- 511
- 512
- 513
- 514
- 515
- 516 Fogg, B. J. *Persuasive Technology: Using Computers to Change What We Think and Do*. Morgan Kaufmann, San Francisco, 2003.
- 517
- 518
- 519
- 520 Habernal, I. and Gurevych, I. Which Argument Is More Convincing? Analyzing and Predicting Convincingness of Web Arguments Using Bidirectional LSTM. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 1589–1599, 2016.
- 521
- 522
- 523
- 524
- 525
- 526
- 527 He, B., Liu, Y., Zhang, Z., Jia, Z., Wu, H., He, Z., Zheng, Z., and Kang, Y. Make an Offer They Can't Refuse: Grounding Bayesian Persuasion in Real-World Dialogues Without Pre-Commitment, October 2025. Preprint at <http://arxiv.org/abs/2510.13387>.
- 528
- 529
- 530
- 531
- 532 Hong, J., Dragan, A., and Levine, S. Planning Without Search: Refining Frontier LLMs with Offline Goal-Conditioned RL, 2025. Preprint at <https://arxiv.org/abs/2505.18098>.
- 533
- 534
- 535
- 536
- 537 Idziejczak, M., Korzavatykh, V., Stawicki, M., Chmutov, A., Korcz, M., Bładek, I., and Brzezinski, D. Among Them: A Game-Based Framework for Assessing Persuasion Capabilities of LLMs. In *Proceedings of the 17th International Conference on Agents and Artificial Intelligence*, volume 15874, pp. 183–195, Singapore, 2025. Springer. doi: 10.1007/978-981-96-8186-0.15.
- 538
- 539
- 540
- 541
- 542
- 543
- 544 Jaipersaud, B., Krueger, D., and Lubana, E. S. How Do LLMs Persuade? Linear Probes Can Uncover Persuasion Dynamics in Multi-Turn Conversations, August 2025. Preprint at <http://arxiv.org/abs/2508.05625>.
- 545
- 546
- 547
- 548
- 549
- Jin, C., Ren, K., Kong, L., Wang, X., Song, R., and Chen, H. Persuading Across Diverse Domains: A Dataset and Persuasion Large Language Model. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 1678–1706, Bangkok, Thailand, 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.acl-long.92. URL <https://aclanthology.org/2024.acl-long.92>.
- Jones, C. R. and Bergen, B. K. Lies, Damned Lies, and Language Statistics: A Comprehensive Review of Risks from Manipulation, Persuasion, and Deception with Large Language Models. *Artificial Intelligence Review*, 59:116, 2026. doi: 10.1007/s10462-026-11517-6.
- Ju, T., Chen, Y., Fei, H., Lee, M.-L., Hsu, W., Cheng, P., Wu, Z., Zhang, Z., and Liu, G. On the Adaptive Psychological Persuasion of Large Language Models, 2025. Preprint at <https://arxiv.org/abs/2506.06800>.
- Kim, T., Lee, J., Yoon, S., Kim, S., and Lee, D. Towards Personalized Conversational Sales Agents: Contextual User Profiling for Strategic Action. In Christodoulopoulos, C., Chakraborty, T., Rose, C., and Peng, V. (eds.), *Findings of the Association for Computational Linguistics: EMNLP 2025*, pp. 5131–5154, Suzhou, China, November 2025. Association for Computational Linguistics. ISBN 979-8-89176-335-7. doi: 10.18653/v1/2025.findings-emnlp.275. URL <https://aclanthology.org/2025.findings-emnlp.275/>.
- Kowal, M., Timm, J., Godbout, J.-F., Costello, T., Arechar, A. A., Pennycook, G., Rand, D., Gleave, A., and Pelrine, K. It's the Thought that Counts: Evaluating the Attempts of Frontier LLMs to Persuade on Harmful Topics, August 2025. Preprint at <http://arxiv.org/abs/2506.02873>.
- Li, Z., Zheng, C., and Shi, Y. Data Speaks, but Who Gives It a Voice? Understanding Persuasive Strategies in Data-Driven News Articles. *IEEE Transactions on Visualization and Computer Graphics*, 2025.
- Liu, M., Xu, Z., Zhang, X., An, H., Qadir, S., Zhang, Q., Wisniewski, P. J., Cho, J.-H., Lee, S. W., Jia, R., and Huang, L. LLM Can Be a Dangerous Persuader: Empirical Study of Persuasion Safety in Large Language Models. In *Proceedings of the Second Conference on Language Modeling*, 2025.
- Ma, Y., Zhang, X., Ren, J., Wang, R., and Wang, M. Linguistic Features of AI Mis/Disinformation and the Detection Limits of LLMs. *Nature Communications*, 16:67145, 2025. doi: 10.1038/s41467-025-67145-1.

- 550 Maier, S., Feuerriegel, S., and Hölbling, L. A Meta-Analysis
551 of the Persuasive Power of Large Language Models. *Sci-*
552 *entific Reports*, 15(1):30783, 2025.
- 553 Meguellati, E., Civelli, S., Han, L., and Bernstein, A. LLM-
554 Generated Ads: From Personalization Parity to Persua-
555 sion Superiority, 2025. Preprint at [https://arxiv.](https://arxiv.org/abs/2512.03373v1)
556 [org/abs/2512.03373v1](https://arxiv.org/abs/2512.03373v1).
- 557 Namikoshi, K., Shamma, D. A., Iliev, R., Fang, J., Filipow-
558 icz, A., Hogan, C. L., Wu, C., and Arechiga, N. Lever-
559 aging Language Models and Bandit Algorithms to Drive
560 Adoption of Battery-Electric Vehicles, 2024. Preprint at
561 <https://arxiv.org/abs/2410.23371>.
- 562 O’Keefe, D. J. *Persuasion: Theory and Research*. Sage
563 Publications, Thousand Oaks, CA, 2015.
- 564 OpenAI. Openai o1 system card. System card, OpenAI,
565 December 2024. URL [https://cdn.openai.com/](https://cdn.openai.com/o1-system-card.pdf)
566 [o1-system-card.pdf](https://cdn.openai.com/o1-system-card.pdf).
- 567 Pauli, A. B., Augenstein, I., and Assent, I. Measuring and
568 Benchmarking Large Language Models’ Capabilities to
569 Generate Persuasive Language. In *Proceedings of the*
570 *2025 Conference of the Nations of the Americas Chapter*
571 *of the Association for Computational Linguistics: Human*
572 *Language Technologies (Volume 1: Long Papers)*, pp.
573 10056–10075, 2025.
- 574 Ramani, G. P., Karande, S., V. S., and Bhatia, Y. Persua-
575 sion Games Using Large Language Models, September
576 2024. Preprint at [http://arxiv.org/abs/2408.](http://arxiv.org/abs/2408.15879)
577 [15879](http://arxiv.org/abs/2408.15879).
- 578 Rescala, P., Ribeiro, M. H., Hu, T., and West, R. Can Lan-
579 guage Models Recognize Convincing Arguments? In
580 *Findings of the Association for Computational Linguis-*
581 *tics: EMNLP 2024*, pp. 8826–8837, 2024.
- 582 Rogiers, A., Noels, S., Buyl, M., and De Bie, T. Persuasion
583 with Large Language Models: A Survey, 2024. Preprint
584 at <https://arxiv.org/abs/2411.06837v1>.
- 585 Salvi, F., Horta Ribeiro, M., Gallotti, R., and West, R. On
586 the Conversational Persuasiveness of GPT-4. *Nature*
587 *Human Behaviour*, 9(8):1645–1653, 2025. doi: 10.1038/
588 s41562-025-02194-6.
- 589 Simpson, E. and Gurevych, I. Finding Convincing Ar-
590 guments Using Scalable Bayesian Preference Learning.
591 *Transactions of the Association for Computational Lin-*
592 *guistics*, 6:357–371, 2018.
- 593 Singh, S., Singla, Y. K., SI, H., and Krishnamurthy, B. Mea-
594 suring and Improving Persuasiveness of Large Language
595 Models, October 2024. Preprint at [http://arxiv.](http://arxiv.org/abs/2410.02653)
596 [org/abs/2410.02653](http://arxiv.org/abs/2410.02653).
- 597 Song, Y. and Wang, H. Would You Like to Make a Dona-
598 tion? A Dialogue System to Persuade You to Donate. In
599 Calzolari, N., Kan, M.-Y., Hoste, V., Lenci, A., Sakti, S.,
600 and Xue, N. (eds.), *Proceedings of the 2024 Joint Inter-*
601 *national Conference on Computational Linguistics, Lan-*
602 *guage Resources and Evaluation (LREC-COLING 2024)*,
603 pp. 17707–17717, Torino, Italia, May 2024. ELRA and
604 ICCL. URL [https://aclanthology.org/2024.](https://aclanthology.org/2024.lrec-main.1540/)
[lrec-main.1540/](https://aclanthology.org/2024.lrec-main.1540/).
- Toledo, A., Gretz, S., Cohen-Karlik, E., Friedman, R.,
Venezian, E., Lahav, D., Jacovi, M., Aharonov, R., and
Slonim, N. Automatic Argument Quality Assessment -
New Datasets and Methods. In *Proceedings of the 2019*
Conference on Empirical Methods in Natural Language
Processing and the 9th International Joint Conference
on Natural Language Processing (EMNLP-IJCNLP), pp.
5625–5635, 2019.
- Wang, J., Yao, Z., Li, L., Qian, J., Yang, Z., and Yu, H.
ChatThero: An LLM-Supported Chatbot for Behavior
Change and Therapeutic Support in Addiction Recov-
ery, October 2025. Preprint at [http://arxiv.org/](http://arxiv.org/abs/2508.20996)
[abs/2508.20996](http://arxiv.org/abs/2508.20996).
- Yao, Z., Chafekar, T., Wang, J., Han, S., Ouyang, F., Qian,
J., Li, L., and Yu, H. ChatCLIDS: Simulating Persua-
sive AI Dialogues to Promote Closed-Loop Insulin Adop-
tion in Type 1 Diabetes Care. In *Proceedings of the*
AAAI Conference on Artificial Intelligence, volume 40,
pp. 39539–39547, 2026.
- Yeginbergen, A., Oronoz, M., and Agerri, R. Dynamic
Knowledge Integration for Evidence-Driven Counter-
Argument Generation with Large Language Models. In
Findings of the Association for Computational Linguis-
tics: ACL 2025, pp. 22568–22584, 2025.
- Yu, F., Jiang, L., Huang, S., Wu, Z., and Dai, X. Persua-
siveToM: A Benchmark for Evaluating Machine Theory
of Mind in Persuasive Dialogues, May 2025. Preprint at
<http://arxiv.org/abs/2502.21017>.
- Zhu, K., Du, H., Hong, Z., Yang, X., Guo, S., Wang, D. Z.,
Wang, Z., Qian, C., Tang, X., Ji, H., et al. Multiagent-
bench: Evaluating the Collaboration and Competition of
LLM Agents. In *Proceedings of the 63rd Annual Meeting*
of the Association for Computational Linguistics (Volume
1: Long Papers), pp. 8580–8622, Vienna, Austria, 2025.
Association for Computational Linguistics.

A. Detailed Scope & Definitions

This appendix provides detailed definitions of key terms used throughout this review.

A.1. Persuasion as an Umbrella Term

This paper focuses on the automated evaluation of LLMs’ ability to influence human beliefs, attitudes, and behaviour through text-based communication. To capture this range of phenomena within a single coherent framework, we use **persuasion** as an umbrella term encompassing all forms of communicative influence that operate through cognitive and decision-making processes, including manipulation and deception, where it is used as a technique of persuasive influence. We adopt this framing following Jones & Bergen (2026), in order to have a term more general than rational persuasion and manipulation, but more specific than influence, which also incorporates related concepts such as coercion and exploitation. The term is value-neutral, capturing prosocial and harmful influence within a single coherent category, and reflects how the broader literature is organised (Jones & Bergen, 2026; Bozdag et al., 2026).

More precisely, persuasion is defined following O’Keefe (2015) as “a successful intentional effort at influencing another’s mental state through communication in a circumstance in which the persuadee has some measure of freedom.”

A.2. Rational Persuasion vs Manipulation

The persuasion term in academic studies is further split into **rational persuasion**, characterised by transparent and autonomy-preserving arguments, and **manipulation**, which covertly exploits the persuadee’s decision-making vulnerabilities in ways the target is not fully aware of and therefore cannot consent to it (Jones & Bergen, 2026).

A.3. Deception

A related concept, **deception**, refers broadly to causing someone to hold false beliefs and is treated, following Jones & Bergen (2026) and El-Sayed et al. (2024), as a special case of manipulation. However, not all deception is persuasive in intent – hallucinations and task-driven deceptive AI behaviour, for instance, are not oriented toward shifting beliefs or behaviour (Jones & Bergen, 2026). This study, therefore, includes only deception that is deliberately used as a technique of persuasive influence.

B. Detailed Search Strategy

This appendix provides details of our eligibility criteria and search methodology.

B.1. Eligibility Criteria Clarifications

Given the early and rapidly evolving state of the field, peer-reviewed publications, preprints, and blog posts were all considered eligible to ensure comprehensive coverage of emerging methods (where a paper was initially identified as a preprint but subsequently published in a peer-reviewed venue by the time of writing, the reference was updated to reflect the published version).

To be included, a method was required to produce a concrete measure of persuasive effectiveness (such as a score, rate, or outcome), rather than analysing persuasion phenomena such as persuasive strategy classification or detection (e.g., Yu et al. 2025). This criterion applied regardless of whether persuasion evaluation was the primary contribution or secondary to another objective, such as developing a more persuasive LLM (e.g., Jin et al. 2024). We restricted our analysis to text-based outputs, reflecting the current state of the field in which fully automated evaluation methods for multimodal persuasion remain limited.

B.2. Search Strategy Details

To ensure comprehensive coverage of relevant literature, we employed a multifaceted search strategy combining traditional academic databases, preprint repositories, and LLM-assisted discovery.

We performed a traditional database search across three major scholarly platforms: Scopus, ACL Anthology, and ArXiv on 10 January 2026, using the following query to search through titles, abstracts and keywords:

660 ("persuasion" OR "persuasiveness" OR "persuasive language") AND
 661 ("evaluation" OR "assessment" OR "scoring" OR "benchmark" OR "measuring") AND
 662 ("language model" OR "LLM" OR "AI")

664 After filtering for publication year (≥ 2021) and language (English only), 116 papers remained from Scopus, 209 from
 665 ArXiv, and 28 from ACL Anthology.

667 **B.3. LLM-Assisted Filtering**

668 For all three database sources, an LLM (GPT-5.2) was used to label each paper as “yes” (to be included), “no”, or “maybe”
 669 given its title and abstract. The prompt used was:

671 Your Task: Determine if the paper proposes or uses a 100% automated method for
 672 evaluating LLM persuasiveness.

673

674 Possible Labels:

- 675 * YES: The paper evaluates LLM persuasiveness without human judgement in the
 676 evaluation process (this excludes human participation in validating the
 677 methodology).
- 678 * NO: The paper does not contain an automated LLM persuasiveness evaluation
 679 technique.
- 680 * MAYBE: The title and abstract lack detail to confirm if there is an automated LLM
 681 persuasion evaluation method present.

682 Output: Provide only one word label: YES, NO, or MAYBE.

683

684 Paper to Evaluate:

685 Title: {Title}

686 Abstract: {Abstract}

687 Across all three sources, 110 were labelled “yes”, 145 “no”, and 98 “maybe”.

688 To verify label reliability, a human reviewer manually inspected a stratified sample of 15 papers (five per label category).
 689 Although the LLM incorrectly labelled several irrelevant papers as relevant, none of the papers it labelled as “no” were
 690 found to be relevant upon manual inspection, providing limited but encouraging evidence that this LLM-assisted labelling
 691 process is effective for exclusion.

692 A final manual screening (of papers labelled “yes” and “maybe”) produced 9 papers from Scopus, 20 from ArXiv, and 3
 693 from ACL Anthology. After deduplication, we were left with 22 papers.

694 **B.4. LLM-Assisted Search**

695 To further reduce the chance of missing relevant work, we conducted an LLM-assisted search using GPT-5.2 and Google
 696 Gemini 3 Pro in Deep Research mode. The prompt used was:

701 Your Task: Identify peer-reviewed research, pre-prints, or technical frameworks in
 702 any other format that propose fully automated methods for evaluating LLM and
 703 GPAI persuasive capabilities. Provide me with a list of at least 20 papers, if
 704 they meet the criteria.

705

706 Inclusion Criteria: The method must be 100% automated.

707

708 Exclusion Criteria: Exclude any methods that require human judgement in the
 709 evaluation stage (validation is accepted).

710

711 Output Format: Provide a table with these columns: Paper Title; Short Evaluation
 712 Method Description; Link

713 This search resulted in 4 additional papers.

B.5. Backwards Reference Searching

Finally, we performed backwards reference searching on a subset of key papers selected as the most comprehensive existing surveys on LLM persuasion (Bozdag et al., 2026; Rogiers et al., 2024; Bassi et al., 2024), which yielded 1 further paper.

B.6. Final Corpus

In total, we identified **27 papers** within the scope, with 3 papers presenting multiple evaluation methods, resulting in **30 distinct evaluation methods**.

C. Persuasion Evaluation Resource Hub

Figure 2 shows an anonymised preview of the online resource hub. The website allows filtering studies by available resources (datasets, codebases) or by evaluation features (personalised persuasion, persuasion strategies).

Figure 2. Anonymised preview of the persuasion evaluation resource hub.

The screenshot shows the 'Resources' page of the LLM Persuasiveness Evaluation Hub. At the top, there is a navigation bar with 'Home', 'Resources', and 'Contact' links, along with search and moon icons. A warning message states: 'This resource list is still being updated and expanded.' Below this, there is a search bar and filter buttons for 'Personalisation', 'Strategy', 'Has data', and 'Has code'. The 'Has data' filter is currently selected. The main content is a table with the following columns: Study, Personalisation, Strategy, Data, and Code.

Study	Personalisation	Strategy	Data	Code
A Large-scale Dataset for Argument Quality Ranking: Construction and Analysis 2020 · Gretz et al.	—	—	IBM-ArgQ	
Among Them: A Game-Based Framework for Assessing Persuasion Capabilities of LLMs 2025 · Idziejczak et al.	—	✓	Game logs	
Can Language Models Recognize Convincing Arguments? 2024 · Rescala et al.	✓	—	Dataset	
Measuring and Benchmarking Large Language Models' Capabilities to Generate Persuasive Language 2024 · Pauli et al.	—	—	Persuasive-Pairs	
Measuring and Improving Persuasiveness of Large Language Models 2024 · Singh et al.	—	—	Dataset	
Measuring the Persuasiveness of Language Models 2024 · Durmus et al.	—	—	Dataset	