

MPhil Politics, Comparative Government

Edward Anders

May 18, 2025

Abstract

Contents

Abstract	i
List of Tables	iii
List of Figures	iv
1 Introduction	1
2 Literature Review	1
3 Theoretical Framework	1
3.1 Theory and Argumentation	1
3.2 Hypotheses	1
4 Case Selection and Data Gathering	2
4.1 Outcome Measures	2
5 Data analysis	4
5.1 Regression Specification	4
5.2 AI-Generated Content Treatment	5
5.2.1 Thermometer Analysis	5
5.2.2 Ordinal Affective Polarisation Analysis	8
5.3 AI-Labelled Content Treatment	11
5.3.1 Thermometer Analysis	11
5.3.2 Ordinal Affective Polarisation Analysis	11
5.4 Additional Analysis	11
5.4.1 Causal Acyclic Testing	11
5.4.2 Agentic-based Modelling	11
Appendix	13
Bibliography	23

List of Tables

List of Tables

1 AI-Generated Content: Thermometer Gap Results 7

2 AI-Labelled Content: Thermometer Gap Results 10

3 (#tab:codebook-table)YouGov UniOM Survey Codebook 13

4 (#tab:ai-balance)Balance Table of Covariates by AI Treatment Group 18

5 (#tab:ai-balance)Balance Table of Covariates by Label Treatment Group 19

6 AI-Generated Content: Thermometer (mostlikely) Results 20

7 AI-Generated Content: Thermometer (leastlikely) Results 21

8 AI-Labelled Content: Thermometer (mostlikely) Results 22

9 AI-Labelled Content: Thermometer (leastlikely) Results 23

10 AI-Generated Content: Agree Out-Party Respect Beliefs 24

11 AI-Generated Content: Trust in Out-Party to Do What Is Right 25

12 AI-Generated Content: Comfort with Child Marrying Opposing Partisan 26

List of Figures

List of Figures

1	Average in- and out-party thermometer net-difference scores	5
2	Thermometer Score Patchwork Plot for AI-Generated Content	8
3	Ordinal Affective Polarisation Patchwork Plot for AI-Generated Content	9
4	Thermometer Score Patchwork Plot for AI-Labelled Content	11

1 Introduction

2 Literature Review

3 Theoretical Framework

3.1 Theory and Argumentation

- Develop the initial arguments and theoretical framework of your project.
- Discuss how your project relates to existing theoretical approaches in the literature and how these are further developed and/or applied in your research.
- This theoretical framework will most likely only be at the preliminary stage, but it is important to outline the relationships between the key actors or variables in your project.
- This ‘model’ does not have to be formal and explicit, but you may find it helpful to specify causal relationships in terms of dependent variable(s) (the outcome) and independent variables (the explanatory factors).
- Formulate some preliminary testable hypotheses derived from this ‘model’.
- Outline the key assumptions of your argument, as well as the limits to your topic (temporally and spatially).

3.2 Hypotheses

- The hypotheses provided should state the relationship(s) expected to be observed between variables being used
- The formulation of a hypothesis should make clear whether it involves one- or two-tailed tests (i.e. predict an increase, decrease, or change in the outcome variable)
- There are two types of hypotheses to consider:
 - **Confirmatory Hypotheses**
 - * The main focus of the study → what the study is designed to test
 - * There should be well-powered analyses of this hypothesis
 - * Should be backed up by strong theory leading to hypotheses *a priori*
 - **Exploratory Hypotheses**
 - * Hypotheses wish to test but are not the main focus of the study
 - * Often seen as secondary hypotheses looking at mechanisms, subgroups, heterogeneous effects, or downstream outcomes

- Should include as many hypotheses as relate to your theory or intervention
- With more than one hypothesis you will need to specify a procedure for handling multiple hypotheses in the inference criteria section of your PAP

Use the paper *Affect, Not Ideology: A Social Identity Perspective on Polarization* to look at theory. As is *The Origins and Consequences of Affective Polarization in the United States*. This has a lot of critiques of the measures of polarization too.

4 Case Selection and Data Gathering

4.1 Outcome Measures

The measures required to understand AI's affect on affective polarisation are multi-faceted. Different measures can be used to understand the primary outcome of affective polarisation; however, the implication of each measure differs. Druckman and Levendusky (2019) clearly outline the best practices for these affective polarisation measures, and how the measures interact. Therefore, this research chooses to follow these measurement recommendations for use in survey self-reporting (Iyengar *et al.*, 2019).

The most common measure of someone's identification with a political party is through a feeling thermometer score. This aims to understand how warmly or coldly someone feels towards the political parties they most and least prefer. The thermometer scores are measured on a scale of 0 to 100, where 0 is the coldest and 100 is the warmest.¹ This survey experiment firstly asks respondents to identify their most and least preferred party (`mostlikely` and `leastlikely`), allowing for in- and out-party identities to be exposed. We then ask respondents to firstly rate how warmly they feel towards each of these party's leaders, `MLthermo_XY` and `LLthermo_XY`, where XY is replaced by each party leader's initials. The use of party-leader thermometers is a common measure, leaning on valence theory's emphasis on the importance of party leaders in shaping party identification and voting behaviour (Garzia, Ferreira da Silva and Maye, 2023).² Moreover, Druckman and Levendusky's (2019: 119) findings show that respondents are more negative towards party elites rather than party voters; thus, the focus on party leaders here helps elicit the more visceral feelings. Alongside these in- and out-group measures, a net-difference score (`thermo_gap`) is also calculated as the difference between the thermometer scores (`MLthermoMean` - `LLthermoMean`) (Iyengar, Sood and Lelkes, 2012).

¹The wording for the thermometer score questions is as follows: "We'd like to get your feelings toward some of our political leaders and other groups who are in the news these days. On the next page, we'll ask you to do that using a 0 to 100 scale that we call a feeling thermometer. Ratings between 50 degrees and 100 degrees mean that you feel favourable and warm toward the person. Ratings between 0 degrees and 50 degrees mean that you don't feel favourable toward the person and that you don't care too much for that person. You would rate the person at the 50-degree mark if you don't feel particularly warm or cold toward the person."

²The Green Party has two co-leaders, Carla Denyer and Adrian Ramsay. Therefore, ratings of both leaders are asked, and the thermometer scores for the Green Party are averaged to create a single score for the party. The variables `MLthermoMean` and `LLthermoMean` are used as the final thermometer measures for in- and out-group thermometer scores.

The next indicator of affective polarisation is a trait-based rating. This measure identifies the traits that respondents associate with opposing parties (Garrett *et al.*, 2014). The limited scope of the survey experiment meant we focussed on the trait of positive trait of *respect*, and whether respondents associated this trait with opposing parties. Respondents were asked: “To what extent do you agree or disagree with the following statement: [leastlikely] party voters respect my political beliefs and opinions.” This question — coded as **agreedisagree** — was asked in a Likert scale format of levels of agreement.³

Additionally, a similar trait-based measure focussed on *trust* was used (Levendusky, 2013). Here, we ask “And how much of the time do you think you can trust [leastlikely] party to do what is right for the country?”. This question was also asked in a Likert scale format, with the options of **Almost never**, **Once in a while**, **About half of the time**, **Most of the time**, and **Always**. This measure is coded as **xtrust**. Along with the thermometer score, the trait-based views of respect, and trust in opposing parties, Druckman and Levendusky (2019: 119) argue that these measures are good, general measures of prejudices held towards opposing parties.

On the other hand, affective polarisation should also be interested in actual tangible discriminatory behaviour. Therefore an emotional, social-distance-based question is included to understand how comfortable respondents are with having opposing partisans in their lives. For example, Iyengar, Sood and Lelkes (2012) popularised the use of the Almond and Verba (1963) five-nation survey question “Suppose you had a child who was getting married. How would you feel if they married a [leastlikely] party voter?”. Coded as **child**, respondents were given options of **Extremely upset**, **Somewhat upset**, **Neither happy nor upset**, **Somewhat happy**, and **Extremely happy**.

- See pre-analysis plan for details of what to include in this section
- Case selection (why focus on the UK?)
 - I am to have external validity to my research/case?
 - Can I make inferences to other cases?
- What is the case?
 - What is the unit of analysis?
 - What is the time period?
 - What is the geographical scope?
 - What are the key variables?
- What data is being collected?
- How is the data being collected?
 - What is the sampling strategy?

³A full breakdown of the survey experiment variables and values can be found in the codebook [Section 5.5](#) in the appendix.

- Note the UK weighting
- Plan for using agentic modelling
 - Why would I use agentic modelling?
 - What is the agentic modelling?
 - How will I use agentic modelling?
- Note the use of a between-subjects survey experiment
 - Deliberately chosen to avoid sensitivity issues noted by Levendusky and Stecula (2021)

5 Data analysis

The following data analyses focus on all outcome measures of affective polarisation to give a holistic understanding of both general and tangible prejudices, and discriminatory behaviours towards the opposing out-group to that of the respondent’s identified in-group. The analysis is split by the treatments being tested: AI-generated content and AI-labelled content. Each treatment is analysed across the outcome measures of thermometer scores, trait-based measures, and social-distance measures.

5.1 Regression Specification

To test the causal Average Treatment Effect (ATE) of respondents being exposed to AI-generated and AI-labelled content on the set of affective polarisation measures, a series of regression models are estimated. The model specification is given by Equation (1):

$$Y_i = \beta_0 + \beta_1 D_i + \beta_2 \mathbf{X}_i + \beta_3 (D_i \times \mathbf{Z}_i) + \varepsilon_i \quad (1)$$

where:

- Y_i takes the outcome variables (`thermo_gap`, `MLthermoMean`, `LLthermoMean`, `agreedisagree`, `xtrust`, and `child`)
- D_i is the treatment recieved (`ai_treatment` or `label_treatment`)
- \mathbf{X}_i is a vector of covariates (see Balance Check in [Section 5.7](#) for details)
- \mathbf{Z}_i is a vector of possible interaction terms between the treatment and moderators
- ε_i is the error term

In this full specification, β_1 estimates the average treatment effect when the moderator(s) are at their reference level. Estimates are calculated with survey-weighted least squares and ordinal logistic models so

results can be generalised to the UK more broadly. β_2 measures the effect of a one-unit change of a covariate on the outcome variable. β_3 captures the treatment effect heterogeneity across different sub-groups of the moderator, where statistically significant non-zero values suggest the ATE is different for different sub-group characteristics.

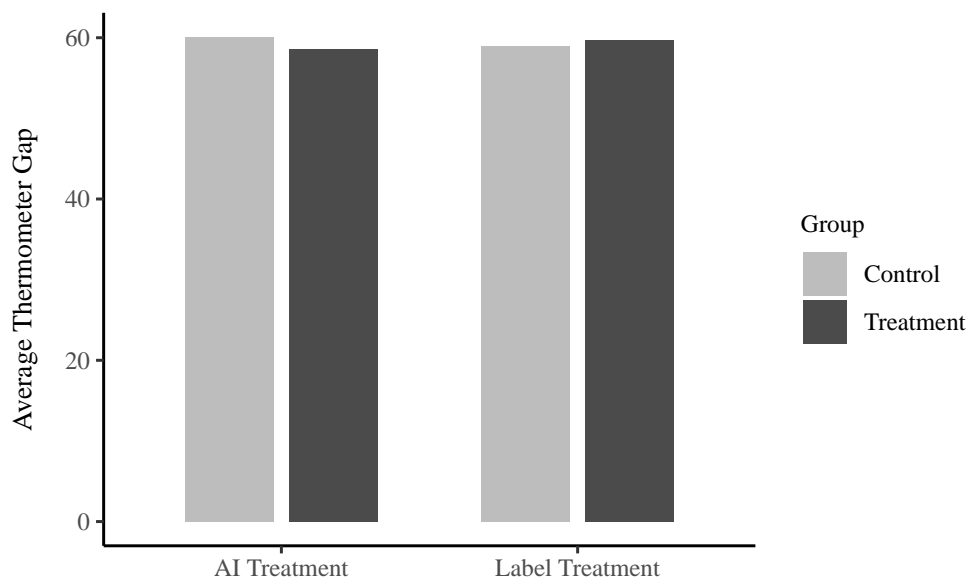
5.2 AI-Generated Content Treatment

The results show no statistically significant treatment effect of AI-generated content on in- and out-party, nor net-difference thermometer scores. However, it is found that Liberal Democrat voters are significantly susceptible to being polarised from exposure to AI-generated content. Trait-based ratings of opposing parties/voters' respect and trust show...

5.2.1 Thermometer Analysis

Thermometer analysis is one of the primary affective polarisation measures. Before determining a causal link between AI content exposure and the affective polarisation measures, a descriptive summary of the `thermo_gap` measures, averaged over all in- and out-party leaders, given for both treatments is presented in [Figure 1](#). This shows how net-difference thermometer scores are similar across both control and treatment groups, suggesting causal effects are likely to be minimal.

Figure 1: Average in- and out-party thermometer net-difference scores



To test whether this descriptive expectation is causally salient, models for the outcome variables for in- and out-party, and net-difference thermometer scores are estimated. The thermometer outcome scores are

continuous measures. Therefore, survey-weighted least squares regression models are estimated.

ATE models are presented in [Table 1](#) for the outcome `thermo_gap`.⁴ A first model (1) sets the benchmark without control for covariates and moderators. A full balance check ([Section 5.7](#)) shows that the treatment and control groups were balanced across all covariates. Despite this, model (2) still includes a full set of pre-treatment covariates as each has theoretical justification for affecting the outcome independently of the treatment, and also to ensure the ATE estimates are efficient. To avoid multicollinearity, individual moderators were sequentially tested within the models; however, few showed any moderation effects. The moderators of party affiliation/warmth (`mostlikely`) and attentiveness to politics (`political_attention`) showed the greatest moderation effects, thus are included in the final model (3) as interaction terms to test these groups for heterogeneity.

The primary takeaway from model (3) in [Table 1](#) is that there is no significant treatment effect seen for the exposure to AI-generated content on the reported net-difference thermometer scores. The treatment group shows a slight decrease of **-0.955 points** in the thermometer gap, but this is not statistically significant. However, respondents who were most likely to vote for the Liberal Democrats showed a significant increase of **11.749 points** at the 95% confidence level in their net-difference thermometer gap, compared to those who were most likely to vote for the Conservative Party. This implies that Liberal Democrat voters are more susceptible to being polarised against their opposing partisans when exposed to AI-generated content. The effect size is notable too. The total treatment effect of AI exposure for Liberal Democrat voters on their affective polarisation is $-0.955 + 11.749 = 10.794$ **points**. This is a significant effect size, and suggests that AI-generated content can have a polarising effect. Other than Liberal Democrat voters, the full model suggests that the treatment of exposure to AI-generated content does not have a significant differential effect on the thermometer gap for different sub-groups of respondents.

To test what drives this affective polarisation, the models for the outcome variables of `MLthermoMean` and `LLthermoMean` are estimated. These models are presented in the Appendix in [Table 6](#) and [Table 7](#). Thermometer scores for both in- and out-party leaders show positive increases of **8.987 points** and **3.13 points** respectively to the treatment but the scores are not statistically significant. These models do not show any sub-group differences from moderator variables.

⁴The full models for the outcome variables of `MLthermoMean` and `LLthermoMean` are available in the appendix in [Table 6](#) and [Table 7](#).

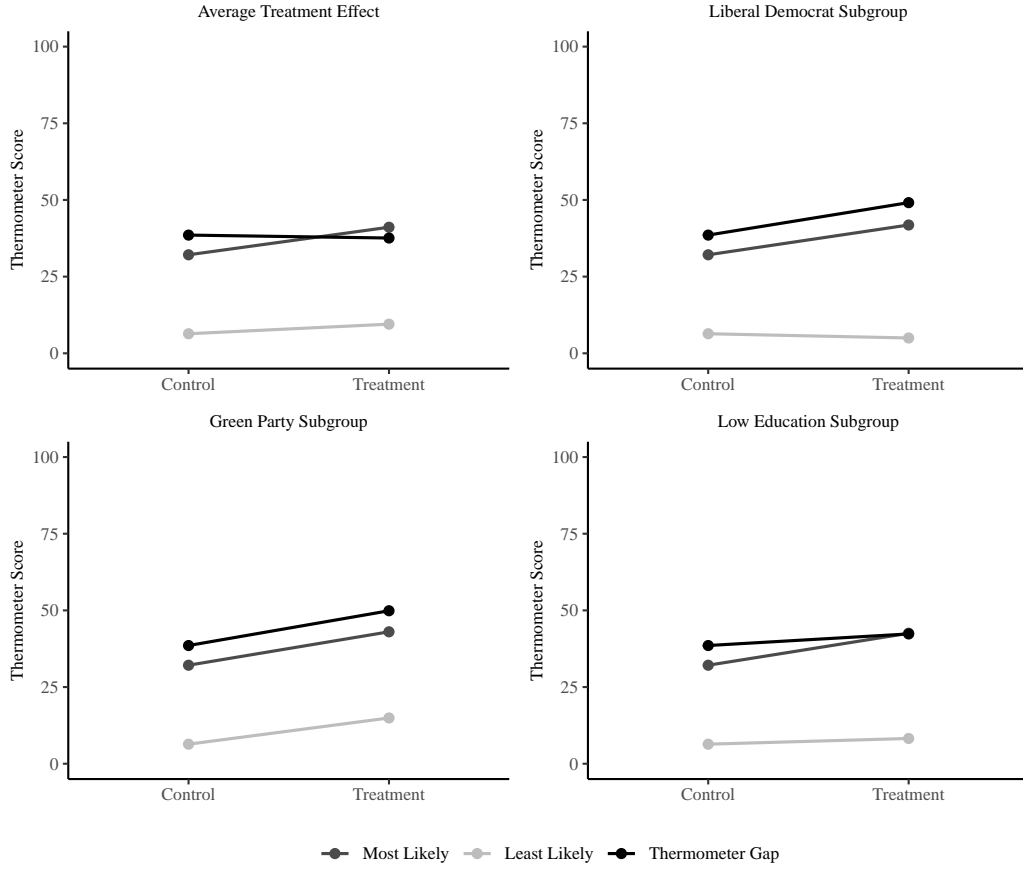
Table 1: AI-Generated Content: Thermometer Gap Results

	Treatment Only	Treatment + Covariates	Full Model
(Intercept)	59.757*** (1.323)	42.045*** (6.399)	38.540*** (7.181)
AI Treatment	-0.491 (1.864)	-0.958 (1.738)	-0.955 (7.701)
mostlikelyGreen Party			15.738** (5.962)
mostlikelyLabour Party			14.057** (4.878)
mostlikelyLiberal Democrats			11.749* (5.250)
mostlikelyReform UK			12.128** (4.054)
AI Treatment:mostlikelyGreen Party			12.293+ (7.312)
AI Treatment:mostlikelyLabour Party			5.179 (4.979)
AI Treatment:mostlikelyLiberal Democrats			11.539* (5.790)
AI Treatment:mostlikelyReform UK			6.251 (5.191)
AI Treatment:Political Attention			-0.719 (0.913)
AI Treatment:EducationLow			4.735 (5.492)
AI Treatment:EducationMedium			-2.912 (3.742)
Num.Obs.	1095	966	966
R2	0.000	0.148	0.189
RMSE	28.31	25.62	25.04
Model	(1)	(2)	(3)

+ p < 0.1, * p < 0.05, ** p < 0.01, *** p < 0.001

Note: Models weighted using YouGov survey weights. The coefficients are reported with robust standard errors in parentheses. Main effects of the included moderators are also reported as rows above the moderator treatment effects.

Figure 2: Thermometer Score Patchwork Plot for AI-Generated Content



5.2.2 Ordinal Affective Polarisation Analysis

To further dissect the affective polarisation caused by treatment to AI-generated content, models for the outcome variables of `agreedisagree`, `xtrust`, and `child` are next estimated. These models are ordinal logistic regression models, as the outcome variables are ordinal measures. The results of these models for each measure of respect, trust, and social-distance to opposing partisans are presented in [Table 10](#), [Table 11](#), and [Table 12](#) respectively.

Figure 3: Ordinal Affective Polarisation Patchwork Plot for AI-Generated Content

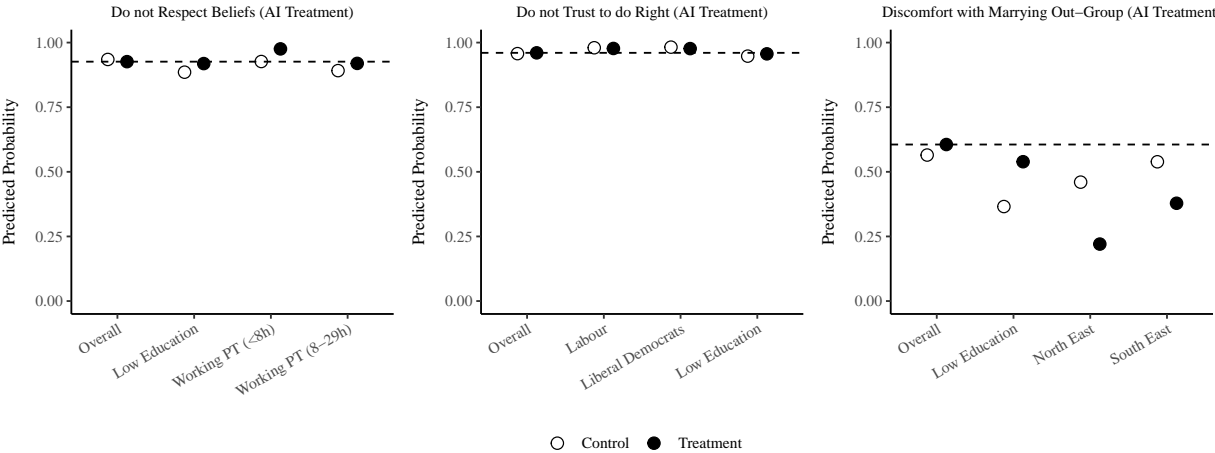


Table 2: AI-Labelled Content: Thermometer Gap Results

	Treatment Only	Treatment + Covariates	Full Model
(Intercept)	59.820*** (1.297)	41.697*** (6.426)	41.616*** (7.132)
Label Treatment	-0.612 (1.862)	-0.540 (1.737)	-7.531 (6.739)
mostlikelyGreen Party			21.545** (7.160)
mostlikelyLabour Party			19.567*** (4.611)
mostlikelyLiberal Democrats			21.799*** (4.881)
mostlikelyReform UK			17.826*** (4.368)
Label Treatment:mostlikelyGreen Party			2.715 (7.571)
Label Treatment:mostlikelyLabour Party			-4.904 (4.908)
Label Treatment:mostlikelyLiberal Democrats			-7.369 (5.766)
Label Treatment:mostlikelyReform UK			-4.386 (5.253)
Label Treatment:Political Attention			1.501+ (0.876)
Num.Obs.	1095	966	966
R2	0.000	0.148	0.187
RMSE	28.32	25.64	25.09
Model	(1)	(2)	(3)

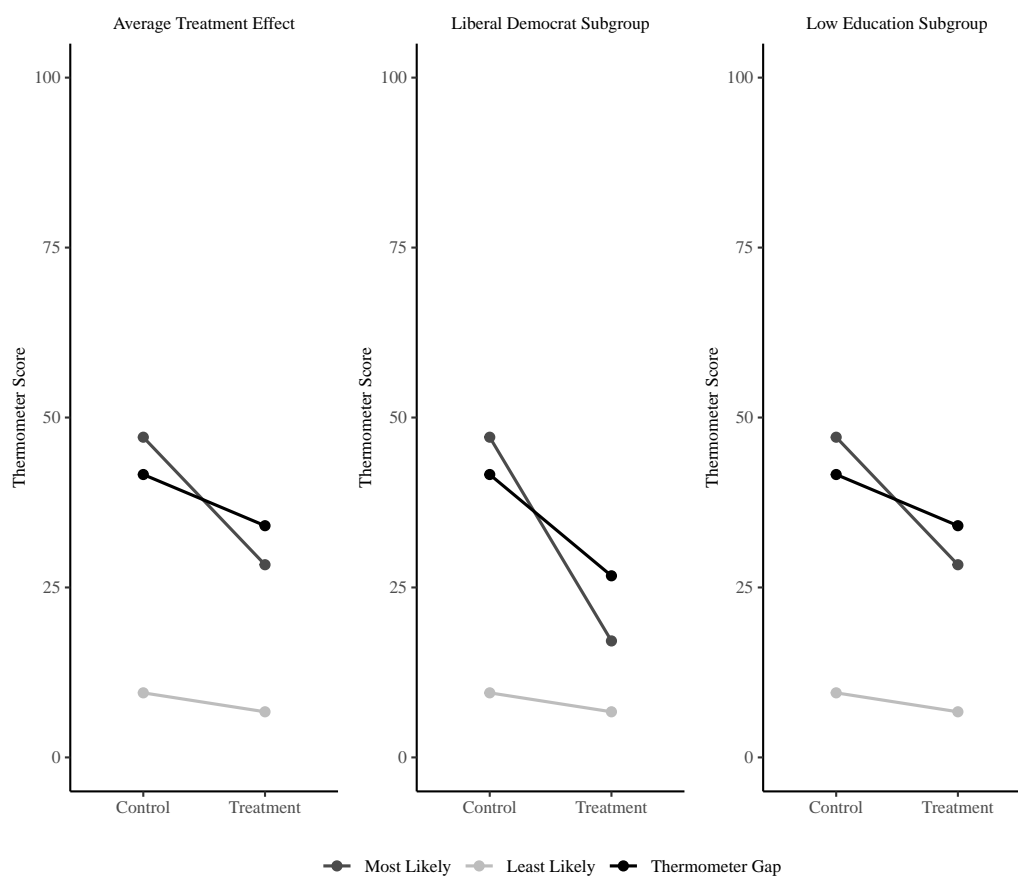
+ p < 0.1, * p < 0.05, ** p < 0.01, *** p < 0.001

Note: Models weighted using YouGov survey weights. The coefficients are reported with robust standard errors in parentheses. Main effects of the included moderators are also reported as rows above the moderator treatment effects.

5.3 AI-Labelled Content Treatment

5.3.1 Thermometer Analysis

Figure 4: Thermometer Score Patchwork Plot for AI-Labelled Content



5.3.2 Ordinal Affective Polarisation Analysis

5.4 Additional Analysis

5.4.1 Causal Acyclic Testing

(see experimental analysis week 6 notes for details)

5.4.2 Agentic-based Modelling

- What is agentic-based modelling?
- Why is agentic-based modelling important?

- How will I use agentic-based modelling?

Appendix

5.5 Codebook

The codebook in Table ?? below provides a summary of the variables used in the YouGov UniOM analysis. The variable names are provided in the first column, followed by the type of variable (e.g., categorical, continuous), a description of the variable, and the values that the variable can take. Note that the outcome variables of `agreedisagree`, `xtrust`, and `child` are ordinal variables on an ordered Likert scale.

Table 3: (#tab:codebook-table)YouGov UniOM Survey Codebook

Variable	Type	Description	Values
<code>identity_client</code>	Identifier	Unique identifier for the respondent	Alphanumeric string
<code>weight</code>	Continuous	Survey weight to ensure national representativeness	Continuous float (e.g., 0.982, 1.034)
<code>age</code>	Continuous	Age of the respondent	Integer values, typically 18–90
<code>profile_gender</code>	Categorical	Gender of the respondent	Female; Male
<code>profile_GOR</code>	Categorical	Government Office Region (region of residence)	East Midlands; East of England; London; North East; North West; Scotland; South East; South West; Wales; West Midlands; Yorkshire and the Humber
<code>voted_ge_2024</code>	Categorical	Did the respondent vote in the 2024 General Election?	Don’t know; No, did not vote; Yes, voted
<code>pastvote_ge_2024</code>	Categorical	How the respondent voted in the 2024 General Election	Conservative; Don’t know; Green; Labour; Liberal Democrat; Other; Plaid Cymru; Reform UK; Scottish National Party (SNP); Skipped
<code>pastvote_EURef</code>	Categorical	How the respondent voted in the 2016 EU Referendum	Can’t remember; I did not vote; I voted to Leave; I voted to Remain
<code>education_recode</code>	Categorical	Re-coded education level (grouped)	High; Medium; Low
<code>profile_work_stat</code>	Categorical	Employment status	Full time student; Not working; Other; Retired; Unemployed; Working full time (30+ hrs); Working part time (8–29 hrs); Working part time (<8 hrs)

Table 3: (#tab:codebook-table)YouGov UniOM Survey Codebook (*continued*)

Variable	Type	Description	Values
<code>political_attention</code>	Continuous	How much attention the respondent pays to politics	Scale (e.g., 0–10 or continuous values)
<code>split</code>	Categorical	Randomly assigned treatment group (1–4)	1 = AI-generated, not labelled as AI-generated; 2 = AI-generated and labelled as AI-generated; 3 = Human-generated but labelled as AI-generated; 4 = Human-generated, not labelled as AI-generated
<code>xconsent</code>	Categorical	Consent to participate in the survey	I consent to taking part in this study; I do not wish to continue with this study
<code>mostlikely</code>	Categorical	Which of these parties would you be most likely to vote for?	Conservative Party; Green Party; Labour Party; Liberal Democrats; Reform UK
<code>leastlikely</code>	Categorical	Which of these parties would you be least likely to vote for?	Conservative Party; Green Party; Labour Party; Liberal Democrats; Reform UK; None of these; Not Asked
<code>MLthermo_KB</code>	Continuous	Thermometer rating for Kemi Badenoch (most likely party)	0–100
<code>MLthermo_KS</code>	Continuous	Thermometer rating for Keir Starmer	0–100
<code>MLthermo_NF</code>	Continuous	Thermometer rating for Nigel Farage	0–100
<code>MLthermo_ED</code>	Continuous	Thermometer rating for Ed Davey	0–100
<code>MLthermo_CD</code>	Continuous	Thermometer rating for Carla Denyer	0–100
<code>MLthermo_AR</code>	Continuous	Thermometer rating for Adrian Ramsay	0–100
<code>LLthermo_KB</code>	Continuous	Thermometer rating for Kemi Badenoch (least likely party)	0–100
<code>LLthermo_KS</code>	Continuous	Thermometer rating for Keir Starmer	0–100

Table 3: (#tab:codebook-table)YouGov UniOM Survey Codebook (*continued*)

Variable	Type	Description	Values
LLthermo_NF	Continuous	Thermometer rating for Nigel Farage	0–100
LLthermo_ED	Continuous	Thermometer rating for Ed Davey	0–100
LLthermo_CD	Continuous	Thermometer rating for Carla Denyer	0–100
LLthermo_AR	Continuous	Thermometer rating for Adrian Ramsay	0–100
agreedisagree	Ordinal	Trait-based measure of whether out-groups respect in-group beliefs	Strongly disagree; Tend to disagree; Neither agree nor disagree; Tend to agree; Strongly agree
xtrust	Ordinal	Level of trust in out-group to do what is right	Almost never; Once in a while; About half of the time; Most of the time; Always
child	Ordinal	Social-distance measure of a child marry an out-group voter	Extremely upset; Somewhat upset; Neither happy nor upset; Somewhat happy; Extremely happy
MLthermoMean	Continuous	Average thermometer score for most likely party	0–100 (row mean of MLthermo scores)
LLthermoMean	Continuous	Average thermometer score for least likely party	0–100 (row mean of LLthermo scores)
thermo_gap	Continuous	Difference between MLthermoMean and LLthermoMean	0–100 (MLthermoMean - LLthermoMean)
ai_treatment	Binary	Treatment status for AI-generated content	1 = Treated (shown AI-generated); 0 = Control (shown human-generated)
label_treatment	Binary	Treatment status for AI-labelled content	1 = Treated (labelled as AI-generated); 0 = Control (labelled as human-generated)

5.6 Data Cleaning

2,001 respondents were provided with the survey experiment. Respondents who did not give consent to participate in the survey were removed. Respondents were given the option to skip questions. When skipped, a value of 997 was assigned to the question, which was then recoded to NA, as were **Not asked** values.

The survey was interested in understanding respondents' views towards their most and least preferred party. When asked who the **mostlikely** and **leastlikely** party was, respondents were given the option to select **None of these**. Respondents who selected **None of these** were removed from the sample as they were unable to answer the follow-up questions.

Categorical variables were recoded to be **factors** in R, these were **profile_gender**, **profile_GOR**, **voted_ge_2024**, **pastvote_ge_2024**, **pastvote_EURef**, **profile_education_level**, **education_recode**, **profile_work_stat**, **xconsent**, **mostlikely**, **leastlikely**, **agreedisagree**, **xtrust**, and **child**.

Each of the thermometer variables were recoded to be **numeric** variables: **MLthermo_KB**, **MLthermo_KS**, **MLthermo_NF**, **MLthermo_ED**, **MLthermo_CD**, **MLthermo_AR**, **LLthermo_KB**, **LLthermo_KS**, **LLthermo_NF**, **LLthermo_ED**, **LLthermo_CD**, and **LLthermo_AR**. As the Green Party has two co-leaders, a mean thermometer score is calculated and used for most and least likely party thermometer scores, coded as **MLthermoMean** and **LLthermoMean**.

For treatment effect analysis, respondents were classified into two treatment groups: those shown AI-generated content (**ai_treatment**), identified where the split variable equalled 1 or 2; and those shown AI-labelled content (**label_treatment**), identified where the split variable equalled 2 or 3. Participants in the other split groups were coded as receiving human-generated or unlabelled content. These variables were coded as binary variables, where 1 indicated the treatment group and 0 indicated the control group.

5.7 Balance Check

To ensure that the randomisation process of the treatment allocation was successful, a balance check is conducted to ensure that the treatment and control groups are comparable in every way other than their treatment assignment status. [Table 4](#) and [Table 5](#) below report the balance of the covariates across the treatment groups. The continuous variables of **age** and **political_attention** are reported as means with the standard deviations in parentheses. The remaining categorical variables are reported as a count from the sample, with the proportions in parentheses. If there was a significant difference between the treatment and control groups, this is indicated with a * for $p < 0.05$, ** for $p < 0.01$, and *** for $p < 0.001$. The balance check shows that randomisation was successful across all covariates for both treatment groups as no covariates were significantly different between the treatment and control groups.

Note that the p-values are reported at the variable level, not for each individual category within a categorical variable. For categorical variables (e.g., gender, vote choice), a single p-value is generated using a chi-squared test, which assesses whether the overall distribution of categories differs between treatment and control groups. The individual category rows are displayed for reference, but since the test is run at the variable level, no p-value is reported for each specific level, giving the **NA** values in the tables.

For each of the categorical variables, there is a base reference category. For example, **profile_gender** uses the base reference category **Male** (reported as **Gender (Male)** in the balance tables). This base acts as the comparison group for the other categories, the p-value compares whether the distribution of the other categories is significantly different from the base category.

5.8 Sensitivity Analysis

Given the nature of the results often being reported as null effects, a sensitivity analysis to determine what the smallest true effect that could have detected 80% of the time is calculated.

5.9 MLthermoMean and LLthermoMean Analysis

The models for the outcome variables of **MLthermoMean** and **LLthermoMean** are estimated using the same model specification as for **thermo_gap** in [Table 1](#). These models are presented in [Table 6](#) and [Table 7](#) respectively.

For the **label_treatment** models, the same model specification is used as for the **ai_treatment** models. The models for the outcome variables of **MLthermoMean** and **LLthermoMean** are estimated using the same model specification as for **thermo_gap** in [Table 2](#). These models are presented in [Table 8](#) and [Table 9](#) respectively.

5.10 Ordinal Affective Polarisation Results

The following results presented in [Table 10](#), [Table 11](#), and [Table 12](#) show the log-odds change in the probability of being in a higher level (a higher threshold cut point) of agreement, trust, or comfort respectively.

Table 4: (#tab:ai-balance)Balance Table of Covariates by AI Treatment Group

Variable	Control	Treatment	p-value	Signif.
Age	52.12 (16.74)	51.56 (16.75)	0.521	-
Political attention	6.69 (1.92)	6.61 (1.98)	0.452	-
Gender (male)	374 (50.1)	412 (54.8)	0.080	-
Female	372 (49.9)	340 (45.2)	NA	
Education level (High)	308 (41.3)	308 (41.0)	0.882	-
Low	151 (20.2)	160 (21.3)	NA	
Medium	287 (38.5)	284 (37.8)	NA	
Employment status (Full time student)	31 (4.2)	35 (4.7)	0.759	-
Not working	32 (4.3)	37 (4.9)	NA	
Other	16 (2.1)	13 (1.7)	NA	
Retired	222 (29.8)	210 (27.9)	NA	
Unemployed	12 (1.6)	21 (2.8)	NA	
Working full time (30 or more hours per week)	327 (43.8)	338 (44.9)	NA	
Working part time (8-29 hours a week)	94 (12.6)	87 (11.6)	NA	
Working part time (Less than 8 hours a week)	12 (1.6)	11 (1.5)	NA	
Voted in 2024 General Election (Don't know)	3 (0.4)	1 (0.1)	0.574	-
No, did not vote	97 (13.0)	102 (13.6)	NA	
Yes, voted	646 (86.6)	649 (86.3)	NA	
Vote in 2024 General Election (Conservative)	162 (25.1)	143 (22.0)	0.587	-
Don't know	2 (0.3)	6 (0.9)	NA	
Green	58 (9.0)	51 (7.9)	NA	
Labour	211 (32.7)	245 (37.8)	NA	
Liberal Democrat	90 (13.9)	84 (12.9)	NA	
Other	13 (2.0)	12 (1.8)	NA	
Plaid Cymru	2 (0.3)	2 (0.3)	NA	
Reform UK	98 (15.2)	96 (14.8)	NA	
Scottish National Party (SNP)	9 (1.4)	10 (1.5)	NA	
Skipped	1 (0.2)	0 (0.0)	NA	
Vote in EU Referendum (Can't remember)	125 (17.0)	132 (17.8)	0.669	-
I did not vote	287 (39.0)	273 (36.8)	NA	
I voted to Leave	323 (43.9)	337 (45.4)	NA	
Region (East Midlands)	49 (6.6)	61 (8.1)	0.376	-
I voted to Remain	89 (11.9)	79 (10.5)	NA	
East of England	94 (12.6)	73 (9.7)	NA	
London	34 (4.6)	26 (3.5)	NA	
North East	83 (11.1)	84 (11.2)	NA	
North West	44 (5.9)	64 (8.5)	NA	
Scotland	109 (14.6)	120 (16.0)	NA	
South East	79 (10.6)	70 (9.3)	NA	
South West	31 (4.2)	35 (4.7)	NA	
Wales	62 (8.3)	66 (8.8)	NA	
West Midlands	72 (9.7)	74 (9.8)	NA	

Note: P-values are from t-tests (continuous) or chi-squared tests (categorical) comparing groups. Significance levels: * $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$.

Table 5: (#tab:ai-balance)Balance Table of Covariates by Label Treatment Group

Variable	Control	Treatment	p-value	Signif.
Age	51.84 (16.62)	51.84 (16.88)	0.996	-
Political attention	6.58 (1.94)	6.71 (1.96)	0.200	-
Gender (male)	408 (54.0)	378 (50.9)	0.240	-
Female	347 (46.0)	365 (49.1)	NA	
Education level (High)	321 (42.5)	295 (39.7)	0.542	-
Low	153 (20.3)	158 (21.3)	NA	
Medium	281 (37.2)	290 (39.0)	NA	
Employment status (Full time student)	31 (4.1)	35 (4.7)	0.966	-
Not working	37 (4.9)	32 (4.3)	NA	
Other	16 (2.1)	13 (1.7)	NA	
Retired	213 (28.2)	219 (29.5)	NA	
Unemployed	19 (2.5)	14 (1.9)	NA	
Working full time (30 or more hours per week)	338 (44.8)	327 (44.0)	NA	
Working part time (8-29 hours a week)	90 (11.9)	91 (12.2)	NA	
Working part time (Less than 8 hours a week)	11 (1.5)	12 (1.6)	NA	
Voted in 2024 General Election (Don't know)	2 (0.3)	2 (0.3)	0.154	-
No, did not vote	113 (15.0)	86 (11.6)	NA	
Yes, voted	640 (84.8)	655 (88.2)	NA	
Vote in 2024 General Election (Conservative)	148 (23.1)	157 (24.0)	0.927	-
Don't know	4 (0.6)	4 (0.6)	NA	
Green	55 (8.6)	54 (8.2)	NA	
Labour	233 (36.4)	223 (34.0)	NA	
Liberal Democrat	85 (13.3)	89 (13.6)	NA	
Other	10 (1.6)	15 (2.3)	NA	
Plaid Cymru	2 (0.3)	2 (0.3)	NA	
Reform UK	96 (15.0)	98 (15.0)	NA	
Scottish National Party (SNP)	7 (1.1)	12 (1.8)	NA	
Skipped	0 (0.0)	1 (0.2)	NA	
Vote in EU Referendum (Can't remember)	131 (17.6)	126 (17.2)	0.490	-
I did not vote	272 (36.5)	288 (39.4)	NA	
I voted to Leave	343 (46.0)	317 (43.4)	NA	
Region (East Midlands)	56 (7.4)	54 (7.3)	0.700	-
I voted to Remain	78 (10.3)	90 (12.1)	NA	
East of England	84 (11.1)	83 (11.2)	NA	
London	32 (4.2)	28 (3.8)	NA	
North East	86 (11.4)	81 (10.9)	NA	
North West	57 (7.5)	51 (6.9)	NA	
Scotland	116 (15.4)	113 (15.2)	NA	
South East	80 (10.6)	69 (9.3)	NA	
South West	28 (3.7)	38 (5.1)	NA	
Wales	72 (9.5)	56 (7.5)	NA	
West Midlands	66 (8.7)	80 (10.8)	NA	

Note: P-values are from t-tests (continuous) or chi-squared tests (categorical) comparing groups. Significance levels: * $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$.

Table 6: AI-Generated Content: Thermometer (mostlikely) Results

	Treatment Only	Treatment + Covariates	Full Model
(Intercept)	68.544*** (1.171)	40.286*** (5.960)	32.118*** (6.736)
AI Treatment	-0.698 (1.603)	-0.527 (1.521)	8.987 (6.476)
Age		0.211** (0.076)	0.301*** (0.089)
mostlikelyGreen Party			12.826* (5.278)
mostlikelyLabour Party			12.443** (3.962)
mostlikelyLiberal Democrats			16.593*** (3.898)
mostlikelyReform UK			13.803*** (3.325)
AI Treatment:Age			-0.162+ (0.093)
AI Treatment:EducationLow			1.435 (4.406)
AI Treatment:EducationMedium			-2.823 (3.345)
AI Treatment:mostlikelyGreen Party			1.911 (6.718)
AI Treatment:mostlikelyLabour Party			-1.817 (4.164)
AI Treatment:mostlikelyLiberal Democrats			0.732 (4.569)
AI Treatment:mostlikelyReform UK			2.790 (4.200)
Num.Obs.	1251	1088	1088
R2	0.000	0.149	0.191
RMSE	23.29	21.45	20.93
Model	(1)	(2)	(3)

+ p < 0.1, * p < 0.05, ** p < 0.01, *** p < 0.001

Note: Models weighted using YouGov survey weights. The coefficients are reported with robust standard errors in parentheses. Main effects of the included moderators are also reported as rows above the moderator treatment effects.

Table 7: AI-Generated Content: Thermometer (leastlikely) Results

	Treatment Only	Treatment + Covariates	Full Model
(Intercept)	9.592*** (0.683)	6.372 (4.331)	6.369 (4.262)
AI Treatment	0.428 (1.091)	1.580 (1.107)	3.130 (4.924)
mostlikelyGreen Party			-11.732** (3.987)
mostlikelyLabour Party			-8.797* (4.142)
mostlikelyLiberal Democrats			-4.275 (3.872)
mostlikelyReform UK			-3.420 (2.729)
AI Treatment:EURef VoteI voted to Leave			-3.951 (3.978)
AI Treatment:EURef VoteI voted to Remain			-2.858 (3.423)
AI Treatment:EducationLow			-1.269 (2.967)
AI Treatment:EducationMedium			3.399 (2.722)
AI Treatment:mostlikelyGreen Party			5.418 (4.754)
AI Treatment:mostlikelyLabour Party			-1.814 (3.639)
AI Treatment:mostlikelyLiberal Democrats			-4.507 (4.010)
AI Treatment:mostlikelyReform UK			0.407 (3.518)
Num.Obs.	1295	1121	1121
R2	0.000	0.073	0.105
RMSE	15.89	14.86	14.77
Model	(1)	(2)	(3)

+ p < 0.1, * p < 0.05, ** p < 0.01, *** p < 0.001

Note: Models weighted using YouGov survey weights. The coefficients are reported with robust standard errors in parentheses. Main effects of the included moderators are also reported as rows above the moderator treatment effects.

Table 8: AI-Labelled Content: Thermometer (mostlikely) Results

	Treatment Only	Treatment + Covariates	Full Model
(Intercept)	69.287*** (1.066)	41.276*** (5.914)	47.104*** (7.389)
Label Treatment	-2.201 (1.596)	-2.315 (1.531)	-18.759* (9.477)
age		0.207** (0.076)	0.121 (0.087)
mostlikelyGreen Party			14.233* (5.830)
mostlikelyLabour Party			13.970*** (3.530)
mostlikelyLiberal Democrats			22.290*** (3.847)
mostlikelyReform UK			18.994*** (3.181)
Label Treatment:age			0.136 (0.094)
Label Treatment:Political Attention			2.244* (0.952)
Label Treatment:mostlikelyGreen Party			-0.192 (6.556)
Label Treatment:mostlikelyLabour Party			-5.777 (4.115)
Label Treatment:mostlikelyLiberal Democrats			-11.210* (4.532)
Label Treatment:mostlikelyReform UK			-7.497+ (4.115)
Num.Obs.	1251	1088	1088
R2	0.002	0.151	0.207
RMSE	23.32	21.47	20.85
Model	(1)	(2)	(3)

+ p < 0.1, * p < 0.05, ** p < 0.01, *** p < 0.001

Note: Models weighted using YouGov survey weights. The coefficients are reported with robust standard errors in parentheses. Main effects of the included moderators are also reported as rows above the moderator treatment effects.

Table 9: AI-Labelled Content: Thermometer (leastlikely) Results

	Treatment Only	Treatment + Covariates	Full Model
(Intercept)	9.935*** (0.725)	8.055+ (4.488)	9.499* (4.836)
Label Treatment	-0.247 (1.091)	-0.688 (1.064)	-2.789+ (1.530)
Label Treatment:GenderMale			3.909+ (2.335)
Num.Obs.	1295	1121	1121
R2	0.000	0.071	0.074
RMSE	15.89	14.86	14.85
Model	(1)	(2)	(3)

+ p < 0.1, * p < 0.05, ** p < 0.01, *** p < 0.001

Note: Models weighted using YouGov survey weights. The coefficients are reported with robust standard errors in parentheses. Main effects of the included moderators are also reported as rows above the moderator treatment effects.

Bibliography

Almond, G.A. and Verba, S. (1963) *The Civic Culture: Political Attitudes and Democracy in Five Nations*. Princeton University Press.

Druckman, J.N. and Levendusky, M.S. (2019) ‘What Do We Measure When We Measure Affective Polarization?’, *Public Opinion Quarterly*, 83(1), pp. 114–122.

Garrett, R.K., Gvirsman, S.D., Johnson, B.K., Tsfati, Y., Neo, R. and Dal, A. (2014) ‘Implications of Pro- and Counterattitudinal Information Exposure for Affective Polarization’, *Human Communication Research*, 40(3), pp. 309–332.

Garzia, D., Ferreira da Silva, F. and Maye, S. (2023) ‘Affective Polarization in Comparative and Longitudinal Perspective’, *Public Opinion Quarterly*, 87(1), pp. 219–231.

Iyengar, S., Lelkes, Y., Levendusky, M., Malhotra, N. and Westwood, S.J. (2019) ‘The Origins and Consequences of Affective Polarization in the United States’, *Annual Review of Political Science*, 22(Volume 22, 2019), pp. 129–146.

Iyengar, S., Sood, G. and Lelkes, Y. (2012) ‘Affect, Not Ideology: A Social Identity Perspective on Polarization’, *Public Opinion Quarterly*, 76(3), pp. 405–431.

Table 10: AI-Generated Content: Agree Out-Party Respect Beliefs

	Treatment Only	Treatment + Covariates	Full Model
AI Treatment	−0.137 (0.137)	−0.115 (0.138)	0.769 (0.622)
Strongly disagree Tend to disagree	0.711*** (0.092)	−0.232 (0.438)	0.264 (0.552)
Tend to disagree Neither agree nor disagree	2.061*** (0.116)	1.148** (0.444)	1.662** (0.560)
Neither agree nor disagree Tend to agree	3.760*** (0.229)	2.859*** (0.455)	3.384*** (0.573)
Tend to agree Strongly agree	4.847*** (0.385)	3.949*** (0.563)	4.478*** (0.683)
AI Treatment:Low Education			−0.687* (0.347)
AI Treatment:Medium Education			−0.103 (0.334)
AI Treatment:Not Working			−0.610 (0.889)
AI Treatment:Other			−1.271 (1.231)
AI Treatment:Retired			−0.410 (0.618)
AI Treatment:Unemployed			−0.956 (1.284)
AI Treatment:Working Full Time			−0.598 (0.620)
AI Treatment:Working PT (8–29h)			−1.573* (0.684)
AI Treatment:Working PT (<8h)			−3.201* (1.611)
Num.Obs.	1334	1334	1334
edf	5	17	26
Model	(1)	(2)	(3)

+ p < 0.1, * p < 0.05, ** p < 0.01, *** p < 0.001

Note: Ordered logistic regression with survey weights and robust standard errors in parentheses. Coefficients represent log-odds of agreement that opposing partisans respect political beliefs. Threshold cutpoints are included but have no substantive interpretation.

Table 11: AI-Generated Content: Trust in Out-Party to Do What Is Right

	Treatment Only	Treatment + Covariates	Full Model
AI Treatment	−0.086 (0.140)	−0.032 (0.156)	−0.216 (0.360)
Almost never Once in a while	0.935*** (0.098)	1.068 (0.814)	0.669 (0.862)
Once in a while About half of the time	2.860*** (0.149)	2.994*** (0.818)	2.622** (0.864)
About half of the time Always	4.035*** (0.278)	4.257*** (0.835)	3.891*** (0.877)
Always Most of the time	4.516*** (0.322)	4.811*** (0.845)	4.446*** (0.886)
mostlikelyGreen Party			−1.446** (0.483)
mostlikelyLabour Party			−1.004** (0.371)
mostlikelyLiberal Democrats			−1.319** (0.434)
mostlikelyReform UK			−0.589+ (0.329)
AI Treatment:mostlikelyGreen Party			0.099 (0.568)
AI Treatment:mostlikelyLabour Party			0.171 (0.430)
AI Treatment:mostlikelyLiberal Democrats			0.626 (0.524)
AI Treatment:mostlikelyReform UK			0.095 (0.428)
AI Treatment:EducationLow			−0.199 (0.450)
AI Treatment:EducationMedium			0.115 (0.351)
Num.Obs.	1346	1119	1119
edf	5	38	48
Model	(1)	(2)	(3)

+ p < 0.1, * p < 0.05, ** p < 0.01, *** p < 0.001

Note: Ordered logistic regression with survey weights and robust standard errors in parentheses. Coefficients represent log-odds of trusting that opposing parties will do what is right for the country. Threshold cutpoints are included but have no substantive interpretation.

Table 12: AI-Generated Content: Comfort with Child Marrying Opposing Partisan

	Treatment Only	Treatment + Covariates	Full Model
AI Treatment	0.041 (0.115)	0.245+ (0.132)	−0.599 (0.583)
AI Treatment:Education LevelLow			−1.032* (0.436)
AI Treatment:Education LevelMedium			−0.430 (0.314)
AI Treatment:RegionEast of England			1.193+ (0.616)
AI Treatment:RegionLondon			0.504 (0.631)
AI Treatment:RegionNorth East			1.921* (0.896)
AI Treatment:RegionNorth West			1.015 (0.628)
AI Treatment:RegionScotland			0.531 (0.693)
AI Treatment:RegionSouth East			1.375* (0.596)
AI Treatment:RegionSouth West			1.243+ (0.637)
AI Treatment:RegionWales			1.028 (0.891)
AI Treatment:RegionWest Midlands			0.313 (0.635)
AI Treatment:RegionYorkshire and the Humber			1.116+ (0.655)
AI Treatment:mostlikelyGreen Party			−0.412 (0.475)
AI Treatment:mostlikelyLabour Party			0.410 (0.401)
AI Treatment:mostlikelyLiberal Democrats			0.388 (0.451)
AI Treatment:mostlikelyReform UK			0.631 (0.431)
Num.Obs.	1408	1203	1203
RMSE	2.30	2.26	2.26
Model	(1)	(2)	(3)

+ p < 0.1, * p < 0.05, ** p < 0.01, *** p < 0.001

Note: Ordered logistic regression with survey weights and robust standard errors in parentheses. Coefficients represent log-odds of comfort with a child marrying an opposing party voter. Threshold cutpoints are included but have no substantive interpretation.

Levendusky, M. (2013) *How partisan media polarize America*. Chicago: The University of Chicago Press (Chicago studies in American politics).

Levendusky, M.S. and Stecula, D.A. (2021) '[We Need to Talk: How Cross-Party Dialogue Reduces Affective Polarization](#)', *Elements in Experimental Political Science* [Preprint].