

PoSSUM: A **P**rotocol for **S**urveying **S**ocial-media **U**sters with **M**ultimodal LLMs

Roberto Cerina

`r.cerina@uva.nl`

Institute for Logic, Language and Computation
University of Amsterdam



Abstract

This paper introduces PoSSUM, an open-source protocol for unobtrusive polling of social-media users via multimodal Large Language Models (LLMs). PoSSUM leverages users’ real-time posts, images, and other digital traces to create silicon samples that capture information not present in the LLM’s training data. To obtain representative estimates, PoSSUM employs Multilevel Regression and Post-Stratification (MrP) with structured priors to counteract the observable selection biases of social-media platforms. The protocol is validated during the 2024 U.S. Presidential Election, for which five PoSSUM polls were conducted and published on GitHub and X. In the final poll, fielded October 17–26 with a synthetic sample of 1,054 X users, PoSSUM accurately predicted the outcomes in 50 of 51 states and assigned the Republican candidate a win probability of 0.65. Notably, it also exhibited lower state-level bias than most established pollsters. These results demonstrate PoSSUM’s potential as a fully automated, unobtrusive alternative to traditional survey methods.

Acknowledgements

The development of this protocol was partially funded by the *Talking to Machines* initiative (<https://talkingtomachines.org>) at Nuffield College, University of Oxford. I am grateful to Prof. Raymond Duch for believing in this project and supporting my efforts to realise it.

1 Introduction

This article describes PoSSUM, an open source¹ protocol to poll social-media users unobtrusively using multimodal Large Language Models (LLMs). The protocol seeks to address the skepticism [42] surrounding Artificially Intelligent (AI) polling by establishing a methodology comparable to that used by traditional pollsters who use online panels [63]. Concerns around AI polling are summarised aptly by this anonymous review to a related paper [14]: ‘... *The goal of polling is quite simple: TO. LEARN. FROM. PEOPLE. NOW. I believe this model here does not learn, not from people, and not now*’. Three necessary conditions emerge from this animated critique – to be a credible alternative to random digit dial (rdd) or self-selected online panels, Silicon samples [4] must enable *novel learning* – i.e. must contain more information than the *mould*² on which they are based; must be *human-aligned* – i.e. conditional on the same generating process, they must produce a distribution of responses which matches that of humans; must be *time-sensitive* – i.e. we must be able to learn about changes in preferences and attitudes over time by studying these samples, and these changes should be reflective of true societal dynamics, rather than artifacts of data engineering.

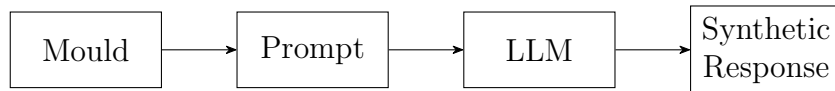


Figure 1: A conceptual description of silicon sampling.

PoSSUM proposes to poll the public by inferring attitudes and preferences of real-life social-media users with multimodal LLMs. A key innovation of PoSSUM entails the use of real-time unstructured digital-trace data to inform a *mould*. The unstructured nature of the data, in its free-flowing and unobtrusively observable nature, provides a much richer compendium of measurable and non-measurable information than the sterile tabular socio-demographic data that has come to define the literature [4, 56, 9]. LLMs can pick up on textual or visual cues that are not easily observable for humans [70]. Importantly, the underlying unstructured data forms a unique digital footprint for a human existing in the world, hence allowing PoSSUM to allocate a unique mould to each individual whose preferences and attitudes we seek to simulate. Others have recently picked up on the potential for informing agents with large unstructured data generated by real human beings (e.g. interview transcripts [50]), but PoSSUM remains alone, at the time of writing, in having these unique moulds update dynamically over time as subjects offer up new text, images, video and other media to the web. This dynamic, time-sensitive mould contributes to solving issues around degradation of simulated agents as we move further away in time from the point at which the mould was created.

A second innovation pertains the generation of representative samples from unrepresentative pools of social media users [2]. Here too PoSSUM leverages the power of LLMs to digest unstructured user profiles and produce socio-demographic labels. This opens the door to

¹<https://github.com/robertocerinaprojects/PoSSUM>

²Silicon samples as per Argyle et al. [4] can be elicited from LLMs by using high-quality, real-life survey responses of humans as a *mould*. See Figure 1 for a graphical illustration of silicon sampling.

traditional non-probability survey sampling techniques for social media users. I introduce a quota sampling module, where quotas are satisfied by matching the users’ inferred characteristics to a stratification frame. Whilst quota sampling can itself produce unrepresentative samples under self-selection [5], it is a valuable tool to partially curb platform selection effects and demographic imbalances [27]. To further address selection and representation issues the protocol relies on Multilevel Regression and Post-Stratification (MrP) [21, 49] with structured priors [19] to analyse the resulting synthetic responses.

PoSSUM’s formulation therefore is set-up to satisfy the three necessary conditions outlined above. It enables *novel learning*, in that the LLMs help us infer previously unknowable tabular survey responses from the unstructured flow of social media data for each user. The temporally updating mould addresses the *time-sensitivity* of learning, clearly enabling the study of preferences changing over time – PoSSUM can learn from synthetic panel data. Finally, efforts to tackle selection and representation via matching individuals to a stratification frame address *human alignment*, under the condition that the LLMs can faithfully infer the distribution of social media users’ preferences from their unstructured data³.

The central argument of this paper is that artificially intelligent polls generated using the PoSSUM protocol are valid instruments for measuring public opinion, on par with traditional polling methods. To test this proposition, PoSSUM is deployed during the 2024 US Presidential election – a setting that allows assessing LLMs’ ability to generalise understandings of political preferences beyond their training data. In this study I employ `gpt-4o-2024-05-13`, whose training concluded in October 2023. Hence the model faces a novel candidate choice-set and unforeseen demographic realignments, characteristic of the 2024 election. Data leakage type criticisms do not apply to this paper as a result.

The paper is organized as follows. Section 2 presents an overview of the PoSSUM protocol. Section 3 describes the `get_pool` routine, which acquires an initial subject pool of social media users. Section 4 introduces the modular prompting architecture employed throughout the protocol. Section 5 reviews the LLM-enabled filtering procedure that selects statistically informative users from the subject pool. Section 6 details the implementation of MrP with structured priors tailored to this framework. Section 7 delineates the evaluation criteria for assessing PoSSUM’s performance, including a multi-dimensional appraisal of predictive accuracy, capacity for novel learning, degree of human alignment, and sensitivity to temporal shifts. Finally, Section 8 provides a comprehensive discussion of key insights arising from the PoSSUM experiment, the limitations of the protocol, and directions for future research.

³LLMs have been shown to be at least as capable as humans on this specific task in previous work [14].

2 Protocol

The goal of PoSSUM is to make granular inference about preferences and attitudes, for a given digital fieldwork period, which is representative of the true underlying population. What I describe is an approach tailored to the \mathbb{X} API, which uses the digital trace of \mathbb{X} users as the mould for LLM generation, but can be extended to any social-media which allows querying of a user panel via user- and content-level queries.

The protocol unfolds as follows: i) *Generating a subject pool* – For each period of interest PoSSUM generates a subject pool of \mathbb{X} users. ii) *Sampling desirable respondents* – Iteratively, each user in the pool is screened for desirable characteristics. Users who do not pass data quality checks, or are unlikely to provide new information to the desired distribution of preferences, are discarded. iii) *Augmenting the mould* – For the surviving users, further information (e.g. their history of tweets) from their \mathbb{X} timeline is elicited. The information is appended to the existing record of the subject, and compiled into a user-specific *mould* – an object containing unstructured multi-media data generated by the user on the social media platform. iv) *Feature extraction* – The mould is passed to a LLM in the form of a prompt. The LLM is given instructions to deduce a set of characteristics for each user according to their mould. The resulting LLM output is structured into tabular data. v) *Hierarchical Bayes* – The synthetic tabular data, representing the hypothetical responses to survey questions from the real-life set of individuals who are active on \mathbb{X} , is then analysed via Hierarchical Bayesian modeling to account for observable selection and address LLM biases. vi) *Post-stratification* – The predictions of the Bayesian multilevel model are post-stratified to generate estimates of the distribution of preferences at the national, state and crosstab levels of analysis.

PoSSUM is therefore composed of three principal routines: `get_pool` (Pseudo-code A.1) is used to identify a subject pool according to keywords matching recent tweets; `poll_users` (Pseudo-code A.2) is designed to implement a series of inclusion checks on the users in the pool, and infer socio-demographics, attitudes and preferences based on their most recent activity on the platform; `make_inference` performs a MrP with structured priors on the synthetic survey data, and generates representative estimates of preferences and attitudes. Figure 2 provides an overview of the protocol⁴.

⁴The protocol is implemented in R, leveraging the `openai` package [55] to call the OpenAI api and prompt the `gpt-4o-2024-05-13` model [47, 46]. Although new versions of the model were released during the course of the study, they were not fit for use to extract user characteristics from social media data. The new models use an instruction hierarchy [66] which flags certain model inputs as high-risk and pushes the model to decline following the instructions. This new model feature generated unacceptable rejection rates, and I therefore opted to keep working with the older model for every PoSSUM poll conducted during the 2024 US election campaign.

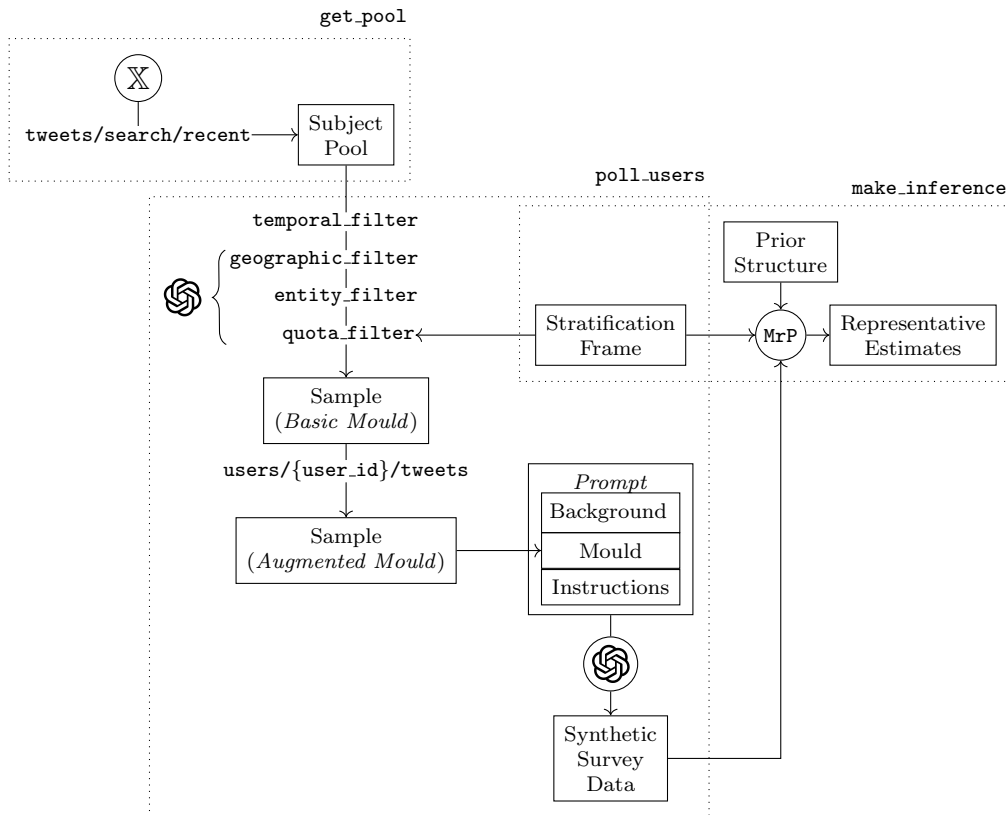


Figure 2: An overview of the PoSSUM protocol.

3 get_pool

I wish to sample a maximally informative set of US adults amongst \mathbb{X} users, to be part of a poll-specific subject pool, from which I can generate a representative sample. Given the time-sensitive nature of the inference I wish to make (what are people’s attitudes **today** ?) I search for these users amongst those who are currently active on the platform.

The \mathbb{X} API *enterprise* and *pro* tiers can be prohibitively expensive, hence I will assume users of PoSSUM have access to the *basic* tier, and cannot collect a simple-random-sample of Tweets from a given day⁵. Even if such sampling were possible, it may not be the default choice. As I demonstrate below, \mathbb{X} API queries can be tuned to target specific sub-populations of active users, potentially mitigating platform-wide selection effects. I propose to use a combination of search queries for the `tweets/search/recent` endpoint⁶, and obtain a series of users who have tweeted the search terms on the platform, up to seven days prior to query-time⁷.

⁵I use the *basic* tier \mathbb{X} api, meaning I pay \$ 100 for downloading 10k tweets per month. I can pay this price multiple times a month, and each payment allows another 10k tweets. The allowance is reset to the original 10k at the end of every month.

⁶See <https://developer.x.com/en/docs/twitter-api/tweets/search/introduction> for more details.

⁷I implement a set of functions, available in the GitHub repository, in the file `X.api.v2.function.R`, reminiscent of the now-defunct `rtweet` [31] and `academictwitterR` [6], in order to specify an appropriate set of queries.

High-Attention Subjects I need the content produced by the selected users to be informative of their political beliefs and attitudes. One way to ensure this is to use political search terms in the \mathbb{X} query. To perform US 2024 pre-election polling we could use a query such as that in Listing A.1. Notice this is a joint query for all the candidates. This is preferable to independent queries per candidate, as these would yield estimates of support subject to selection on the dependent variable. The independent approach ignores the distribution of the search terms across tweets, and over-samples supporters of each candidate, distorting the distribution of support in favour of smaller parties⁸. I assign a *weight* (maximum number of tweets extracted) to this query of size ω .

Low-Attention Subjects Individuals who talk about politics on \mathbb{X} are still unlike their counterparts in the general population. In particular, these are high political attention individuals, who are significantly more likely to vote than their population counterparts. To alleviate selection on political-attention I rely on a second set of queries, which are more likely to sample *normies*. I extract a random sample of L *trending topics* in the US (obtainable via <https://trends24.in/united-states/>), and produce a separate query for each topic. Each trending query is assigned a weight of $\frac{\omega}{L}$, such that the queries seeking to capture high-attention individuals and those capturing *normies* are assigned the same weight. Note this is arbitrary – I noticed this worked well in the US, but in general ω and L are hyper-parameters that need tuning. I end up with a set of queries \mathbf{q} , which is an object of size $K = L + 1$, and a corresponding set of weights $\mathbf{w} = \{\omega, \frac{\omega}{L}, \dots, \frac{\omega}{L}\}$.

I execute each of these queries in a loop, and for each I obtain a tweet-user object $(\mathcal{T}, \mathbf{v})_{kt}$ containing at most w_k tweets, and their associated \mathbb{X} user profiles. The result of `get_pool` is a dated subject pool, containing the profile information about each user (e.g. self-reported description, location, profile picture, etc.), the date on which this user was included in the pool, the search-query used to capture them, and the set of query-related tweets which the user is responsible for.

⁸On the other hand, this sampling is very efficient per party – if you have access to selection-correction terms in the style of King & Zeng [32, 14], this approach would allow for the most sampling-efficient analysis.

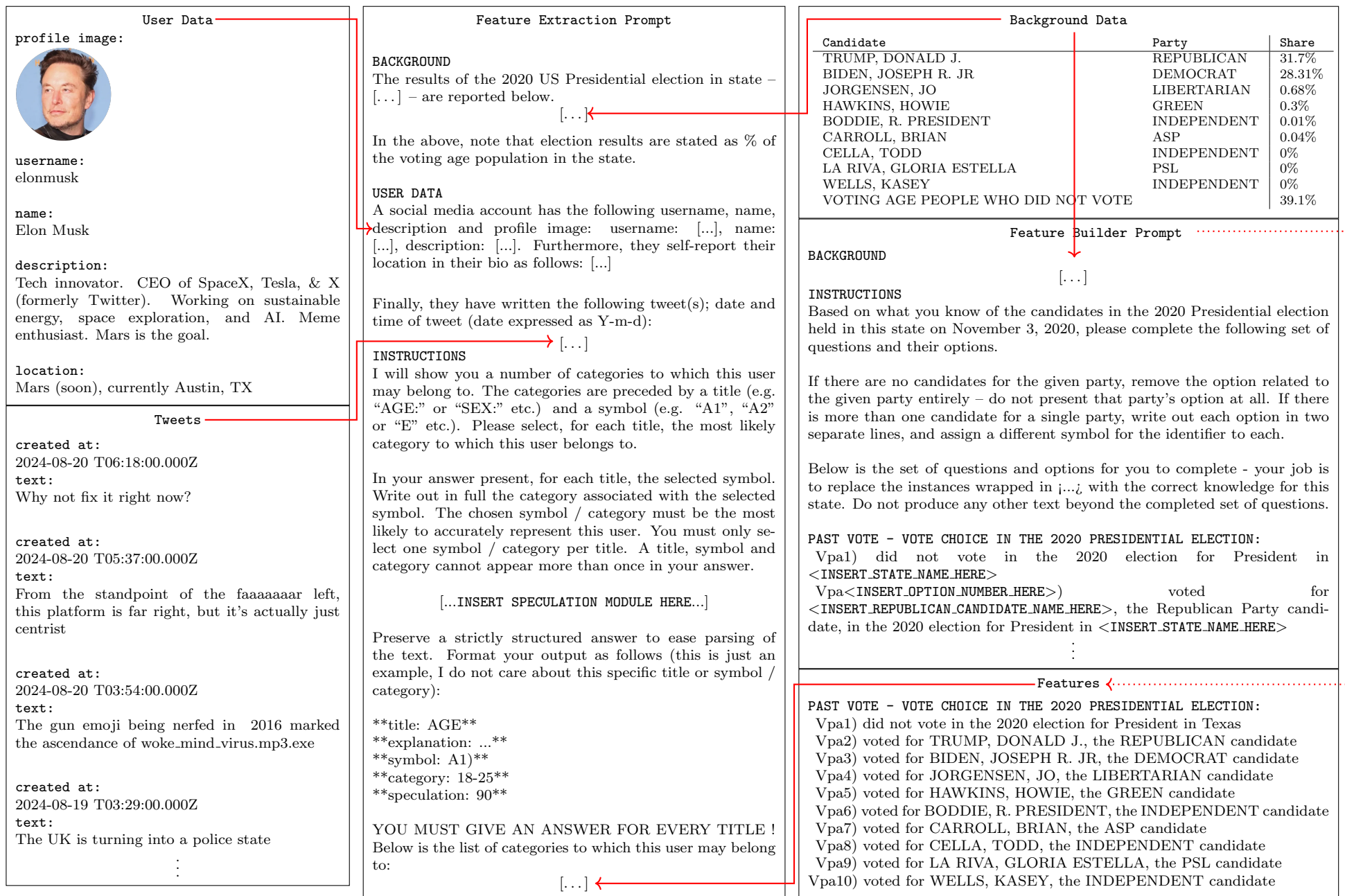


Figure 3: Toy example showing the composition of a prompt under the PoSSUM framework. Red arrows pointing to ‘[...]’ indicate instances where modular prompt components are slotted in. The dotted arrow indicates the LLM generation conditional on a given prompt. The above toy example showcases a single feature (2020 vote choice), though multiple features can generally be extracted simultaneously. Not every prompt contains all of the elements indicated in this Figure.

4 Prompting Architecture

PoSSUM prompts the LLM for a variety of reasons. An initial reason is the filtering of users – non-persons, non-members of the population of interest, and over-sampled users are identified with the aid of LLMs and discarded. I will cover these set of prompts in Section 5.

The principle task of the protocol is *feature extraction*, where the LLM is given a user’s timeline to read and infer a set of socio-demographic characteristics, preferences and attitudes. This is performed in a **feature extraction prompt**. An example of a feature extraction prompt is provided in Figure 3.

The customisation of downstream feature extraction prompts conditional on some user-specific characteristics can also be achieved with a **feature builder prompt**. Figure 3 presents an example of this in reference to the feature “2020 vote choice”. The user in the example lives in Texas, hence the “2020 vote choice” feature-set available to this user should only include candidates which had ballot-access in the state. This is implemented via a two-step prompting strategy: first I generate a prompt that retrieves up-to-date 2020 ballot-access data, and provides a clear example of how to structure a feature-set. The model responds to the prompt by generating a user-specific feature-set; second, the feature set is passed on to the feature extraction prompt (Pseudo-code A.6).

It follows that generally, within PoSSUM, prompts have a standard modular form – they are composed of: i. *background information*; ii. a *mould* based on the available user data; iii. and a set of *instructions*, usually dependent on a set of features of interest. Pseudo algorithm A.3 describes the building of the prompt.

Instructing a Neutral Annotator gpt-4o-2024-05-13 is tasked with annotating unstructured social media profiles of selected users. LLMs are capable general purpose task solvers [10, 1]. Annotation of unstructured data is one such task in which LLMs have shown superhuman performance [24, 61]. Despite some efforts to provide basic best practices around prompting LLMs for this kind of task [60], the degrees of freedom around prompt building are simply too large, and the literature too young, to have had the full spectrum of practices systematically tested. There have been attempts at standardised, systematic testing of prompts [7] though these have largely focused on prompt stability rather than scoring of wholly different architectures. Systematic architecture testing requires the ability prompt the LLM with the same information, but under different architectures, many times over – a prohibitive enterprise when the mould is large. A consequence of this is that the *agentic* approach to silicon sampling [4, 3, 56, 9], where the LLM is asked to role-play, or impersonate a given individual according to a set of characteristics, has not been systematically tested against a simple approach where the AI is a neutral annotator. Moreover the literature is afflicted by baffling findings regarding the sensitivity of prompting strategies to any number of arbitrary tricks. Studies have reported that introducing emotional stimuli in prompts can improve their performance on benchmarks [37], or that treating the AI as a method actor, and providing “dramatic scene settings and role definition” [17] outperforms traditional prompting styles.

I propose a modular prompting architecture, which builds on the Chain of Thoughts (CoT)

[71, 68] approaches to address the social media feature extraction task. I treat the LLM as a *neutral annotator*, purposely avoiding imbuing it with specific personalities. This is similar to a silicon forecaster [58], though the LLM is not explicitly instructed to play the role of an expert, forecaster, or any other.

Examples of relevant instructions are available in Figure 3. Key to the feature-extraction exercise are “feature objects”, which define a specific choice-set. Feature objects are assigned a standard modular structure: each object contains a *title*, which describes a survey question; a set of *categories*, which represent the potential response-set; and each category is identified by a unique *symbol*.

The feature object structure described above is made explicit to the LLM via a set of instructions, in an attempt to pre-emptively address issues around symbolic binding [52]. Strong language is used to encourage the LLM to provide rigidly structured and consistent output, temperature notwithstanding.

An important caveat specific to LLM feature extraction pertains the order in which text is presented in the LLM’s prompts and outputs. The auto-regressive nature of LLMs [35] implies that when text is generated in response to a given prompt, earlier tokens will affect the next-token-probabilities downstream [39]. To encourage the LLM to provide answers which are conditional on some degree of reasoning, rather than the product of post-hoc justification, I enforce an output structure requiring the LLM to provide an explanation before selecting a symbol / category for a given feature. I am limited in the extent of reasoning I can stimulate by the OpenAI API’s limit on output tokens (4,096).

When multiple features are to be extracted simultaneously, the respective feature objects are appended to the prompt. The feature extraction operation then considers all features jointly, and prompts the LLM to produce a joint set of imputed features for the given user. I find for most tasks, simultaneous feature extraction is preferable to a set of independent prompts, one for each attribute of interest. Separating prompts is an intuitively attractive choice due to its preservation of full-independence between extracted features. But this is extremely inefficient in terms of tokens, given that each prompt has to re-describe the background, the mould and the operations of interest. Prompting the LLM to extract all features simultaneously, by including the full list of desired features in a single prompt, is generally a productive approach. To minimise the effects of auto-regression on the generated survey-object, we can randomise the order of all features in the feature extraction prompt, so that order effects on the overall sample cancel-out with a large enough number of observations.

Feature Objects Listing B.3 presents an example of a multi-features object, to be appended to the instructions module of the prompt.

Categories in the features object can seem needlessly verbose – there are two reasons for this: i. detailed descriptions of the categories can help reduce the “neutrality bias” of the LLM – namely the tendency for the LLM to systematically prefer a more “neutral”, “majority class”, or “wide-net” option under uncertainty; ii. associating each category with a unique text string helps ensure the unique parsing of the LLM output, especially when dealing with long prompts which include multiple questions sharing the same answer-text.

It can be helpful at this stage to categorise features of interest in two distinct sets: those

whose distribution in the population is known, and is directly available to us via auxiliary surveys, census data, election results, etc., we consider *independent* variables; those whose distribution in the population is unknown, we consider *dependent* variables. When ascertaining if a given feature belongs to either group, a simple rule of thumb is used: could I weight the poll by the marginal distribution of this variable? if so, this is an independent variable; if not, I consider it dependent. Dependent variables should always come after the independent variables in the prompt, so that their distribution can be conditionally inferred.

Handling Prompting Strategy Uncertainty There is unresolvable uncertainty around LLM queries to infer voting preferences for 2024 (inferred prospective vote) and 2020 (inferred past vote). Small changes in prompt wording can lead to large variations in results [7]. In the context of inferring voting preferences, at least four discernible strategies emerge based on the prompting framework outlined above:

a. *minimally informative* prompting offers a standard choice-set to the LLM, applicable to all users in the pool. The resulting prompt is unconditional of any user characteristic, and vote choice is inferred independently of other feature-extraction tasks. This is desirable if we want the inferred vote to be the result of an exclusive analysis of the mould, and be unaffected by other sources of information;

b. *moderately informative* prompting uses the feature-builder module to update the general choice-set, reflecting candidate options available to the user conditional on some characteristics. In the context of vote choice, the choice-set is conditional on the user’s state of residency;

c. *highly informative* prompting includes the moderately informative provisions, as well as leveraging relevant area-level election results. These past election results are included in the feature extraction prompt to directly influence LLM inference. Conditioning on relevant background election results induces a behaviour similar to raking, in that the inferred individual-level distribution of the vote will be somewhat constrained by the available area-level marginal distributions;

d. *joint socio-demographic* prompts do not use background data, and rely on the choice-set approach of the minimally informative prompting style. The key difference is that vote choice is estimated contemporaneously as the other features of interest, effectively conditioning LLM inference on sequentially inferred socio-demographics. This approach can be useful if the LLM can infer auxiliary characteristics with little error, and if the LLM’s understanding of the relationship between these inferred auxiliary characteristics and the vote is accurate enough to provide useful inferential constraints. Again I wish for the LLM to implicitly perform raking to its internal vote choice representation – the inferred vote distribution is to conform to the underlying marginal distributions of the vote by the inferred socio-demographics, which exist in the LLM’s silicon mind.

In absence of a clear preference amongst prompting strategies, I rely on a wisdom-of-the-synthetic-crowd approach [58, 62, 45]. The heuristic here is that prompt heterogeneity can make LLM learners more uncorrelated, and aggregates of uncorrelated learners typically have desirable properties [25, 28, 29]. I apply a classic majority voting algorithm, breaking ties at random.

5 Filters

ℳ accounts can make for noisy subjects, and not all accounts are deserving of resource allocation for profile augmentation with timeline data. A large number of accounts are not extensions of real existing individual persons, but rather represent organisations [69], bots [18], parody accounts or other non-person entities. A small number of accounts is responsible for a disproportionate amount of activity on the platform [33], and could dominate inference in absence of appropriate mitigation measures. Many accounts contain no discernible information about the location of the user [40], a key feature for being able to make representative inference at the area-level. Relative to the population of interest, select socio-demographic groups are likely to be over-sampled [27, 41].

The `poll_users` routine implements a number of filters to decide which accounts are most deserving of attention – which are most valuable at a given moment in the digital fieldwork period to construct a representative sample of US adults. Some of these filters are *mechanical*, in that they simply apply rules to user meta-data to discard or retain profiles. Others are *intelligent*, in that they leverage AI to deduce some key features of the profile to make a value determination. The filtering routines which follow are applied to the data generated by the `get_pool` routine, so the user-profile data along with usually a single tweet related to the query of interest.

Temporal Filter (Pseudo-code A.4): A routine limiting the number of synthetic survey responses we wish to obtain from a single user within a given time-frame. For example, in the context of pre-election polling, we may wish to collect new data on a given user only once every 30 days – if the digital fieldwork is spread over a full month – despite their more frequent content creation. The routine to implement the temporal exclusion criteria involves: i. tallying the users that have been processed up to now; ii. identifying which of those have been processed within the last 30 days (or whatever the exclusion criteria); iii. removing those users from the *fresh* pool generated by the `get_pool` routine.

Null Geography Filter (Pseudo-code A.5): This is a relatively simple data-quality check. Geography is a fundamental part of pre-election opinion polling – we must be able to place individuals within the given geographic boundary we wish to make inference for. If the user has no self-reported location, we exclude the user a-priori. Intelligent geographic filtering is in principle possible when an explicit location field is absent, by prompting the LLM to infer a location from other content generated by the user – this tends to be less accurate and more expensive due to the larger amounts of input-tokens necessary.

Entity Filter (Listing 1): Consists of defining the kinds of social-media profiles we want to include in our analysis. For pre-election opinion polling, we would wish to exclude ℳ accounts related to organisations (e.g. news outlets, NGOs), bots, and focus solely on real-life persons [69].

Listing 1: Entity Filter prompt.

1	Is this the account of a real-life existing Person, or of another kind of entity ?
2	Respond either with "P" for Person or "0" for Other.

Intelligent Geographic Filter (Listing 2) : This filter helps exclude users who are unlikely to reside in the Level 1 geography. Level 1 constitutes the broadest boundary within which individuals belonging to the population of interest fall. For US pre-election polling I set this to the “United States of America”. Level 2 geography is then intended to be the “State”, and Level 3 is the relevant “Congressional District”, and so on. It is efficient to use a prompt which allows rejection of users who fail the Level 1 inclusion criteria, and simultaneously extracts the Level 2 information. Listing 2 presents an implementation of the geographic extraction prompt. PoSSUM rejects users who are “*Not from a state in the USA*”. The great advantage of using \mathbb{X} relative to other platforms is the relatively high rate of available self-reported locations, which makes geographically-bound polling possible.

Listing 2: Intelligent Geographic Filter prompt.

```

1 Which state of the USA do they live in?
2 If they do not specify a state, but are still from the United States, write "USA".
3 If they are not from a state in the USA, write "Not from a state in the USA".
4 Write out just the full name of the state.
5 If they are from the District of Columbia, also known as Washington D.C., write "District of Columbia".

```

Quota Filter: The population of \mathbb{X} users is notoriously unrepresentative of the US population [27]. It is nevertheless a very large pool of US residents, accounting for around 22% of the US population. Whilst some categories – namely higher educated and higher income individuals – are extraordinarily over-represented, the pool is “deep enough” that we could expect to eventually find a number of representatives for most relevant socio-demographic groups in the population. It follows that implementing quotas is liable to make sampling more efficient.

PoSSUM implements quota sampling as follows: i. define a stratification frame (e.g. Table 1) which describes the number of individuals ω_c^* from each “cell” $c \in \{1, \dots, C\}$, which we would expect to capture in a random sample of target size Ω^* users – we could for instance set $\Omega^* = 1,500$ to produce polls of a somewhat traditional sample size; ii. a feature extraction operation is deployed to infer the values of the relevant variables for the user at hand. At this stage the LLM does not make use of any background information, and it utilises the same user-level information as the other intelligent inclusion criteria; iii. if the user belongs to a cell in the stratification frame for which the number of sampled users ω_c' is smaller than the number of desired users ω_c^* , I retain the user and update the quota counter – otherwise I exclude the user from the analysis. Pseudo-code A.2 contains a symbolic description of the quota exclusion criteria implemented here.

Surviving user profiles are sufficiently information-rich, representing a real-life person in the Level 1 geography. Their latest digital trace is recent, at most 1 week old from the moment the `get_pool` routine is initiated. It is then efficient to expend resources to “survey” these profiles. The LLM is prompted under the general feature extraction framework described in Section 4, with two important differences: i. the digital trace available for each user is expanded further by querying their timeline for their last m tweets, further augmenting their respective mould; ii. we impute a complete set of *independent* and *dependent* characteristics conditional on this new expanded mould.

When expanding profiles, I distinguish between users captured via *trending topics*, as opposed to *political talk*, queries. A very small number of tweets is necessary to estimate the

Cell	Sex	Age	Household Income	Race/Ethnicity	Vote 2020	Quota	Counter
1	male	65 or older	up to 25k	black	D	2	0
2	female	25 to 34	between 25k and 50k	white	D	3	3
3	male	35 to 44	between 75k and 100k	hispanic	D	2	2
4	female	45 to 54	between 75k and 100k	white	D	6	6
5	female	35 to 44	between 25k and 50k	black	D	1	1
\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots
430	female	25 to 34	between 25k and 50k	asian	stayed home	1	0
431	female	65 or older	between 50k and 75k	hispanic	stayed home	1	0
432	female	18 to 24	more than 100k	asian	stayed home	1	0
433	male	18 to 24	between 50k and 75k	native	stayed home	1	0
434	female	55 to 64	between 50k and 75k	asian	stayed home	1	0
435	male	18 to 24	between 50k and 75k	asian	stayed home	1	0

Table 1: Example implementation of a stratification frame with quota counter, for a target sample size $\Omega^* = 1,500$. This is a snapshot taken with 647 respondents still to be collected.

preferences of those who talk politics on \mathbb{X} . Conversely, users discussing trending topics on \mathbb{X} can be totally enigmatic with respect to their politics – their last m tweets could never mention anything remotely useful to indicate political preferences. As a result I set two distinct values of m for these two sets of subjects: $m^{\text{trending}} = \lambda \times m^{\text{politics}}$, $\forall \lambda > 1$. I use $\lambda = 2$ and $m = 20$, but this is open to further tuning. What is generally true is that, where resources permit, “more is better” in terms of information used to generate or expand a user’s mould.

6 make_inference

The goal of PoSSUM is to make representative inference at the population-level, as well as for granular socio-demographic and geographic segments of the population of interest. Each poll will seldom be large enough to make such granular inference, in that crosstabs will be scarcely populated, and not immediately generalisable. Moreover selection effects are still likely to plague the sample due to social media selection, attention selection, and other socio-demographic sample imbalances.

The `make_inference` routine implements a weighting strategy to generalise the findings from silicon samples to the population of interest. The weighting method of choice here is Multilevel Regression with Post-Stratification (MrP) [21, 49, 34]. The explicit knowledge of unfilled quotas prompts a treatment of these cells as having missing dependent variables. We can then use a hierarchical model, under the ignorability assumption [65], to estimate the dependent values for the incomplete cells, and stratify these estimates to obtain national and state-level estimates. This also allows a comprehensive treatment of uncertainty at the cell-level, which is liable to provide more realistic intervals on the poll’s topline than traditional adjustments.

Structured priors [19], as well as deliberate model selection, are crucial – expecting a biased sample and noisy crosstabs, I leverage informative priors and theoretical knowledge of the functional form of the dependent variable to “direct” learning towards a useful configuration. The model is estimated using the Bayesian probabilistic programming language

Stan [11]. Post-stratification is performed at various levels of analysis – predictions from this model are made for every ‘cell’ in the stratification frame, and these are then aggregated at the national, crosstab (e.g. by categories of age, gender, ethnicity, etc.) and state levels respectively, according to the weight associated with each cell (see Cerina & Duch [14] for a comprehensive theoretical model connecting unrepresentative social media samples to post-stratified estimates).

6.1 Extending the Stratification Frame

To improve the MrP estimates I use a modified MrsP [36] procedure (Smooth MrsP)[12, 13]. The goal of this procedure is to extend the stratification frame, which is derived from the 2021 American Community Survey [64], to include the joint distribution of 2020 Vote Choice as derived from an auxiliary survey. It differs from traditional MrsP in that it doesn’t use the auxiliary survey crosstabs to augment the frame, but rather it fits a model to smooth the crosstabs first, and then projects these onto the existing frame. This approach can help generate more plausible estimates for ‘noisy’ cells, when the number of cells in the frame is large and the sample-size-per-cell in the auxiliary survey is relatively small. I use the 2022 Cooperative Election Study (CES) [57] as the auxiliary survey to get estimates of 2020 recall vote⁹. I fit a deep-MrP [23, 26] model using **Stan** [11] to generate estimates of past-vote which leverage interactions between demographics as much as possible, in order to avoid attenuation bias in the estimated cell-level distribution. The likelihood of the model is categorical, and SoftMax is used as the link-function. The *depth* of the Bayesian Hierarchical model is given by the inclusion of marginal effects of sex, age, ethnicity, education, household income and state, as well as all two- and three- way interactions. All effects are estimated as random effects under non-centered parametrisation and recommended weakly-informative priors [20]. The **Stan** code for this model is available in the GitHub repository under name `model_ai.survey.SmoothMrsP.stan`. The resulting frame is then raked to the known state-level distribution of demographics and past vote, using the **anesrake** procedure [51]. The quota-frames used for the quota-filter are samples of size Ω^* from this “mother-frame”, where a new “daughter-frame” is sampled to generate targets for each new poll.

6.2 Hierarchical Bayes to Model Silicon Samples

The final hierarchical model used to generate smoothed estimates of the dependent variable of interest is a simple MrP with structured priors [19]. The “structure” of the model plays an important role here, as it can help smooth the learned effects of a model trained on AI generated data in a sensible way. LLMs can leverage stereotypes in making their imputations [15], which can translate to exaggerated relationships between covariates and dependent variables. Adding structured smoothing to the model allows us to correct for this phenomena, to some degree.

⁹I use this dataset for the following reasons: a. it is a large sample of 60k subjects, affording greater scope for estimating interaction effects between demographic attributes; b. the alternative (ANES) was much too biased in favour of the Democratic candidate in 2020; c. it allows me flexibility to extend the frame further by 2022 vote, using the same dataset, if it is reasonable to do so at a later stage.

<i>predictor</i>	<i>level</i>	<i>description</i>	<i>index</i>	<i>domain</i>	<i>parameter</i>	<i>prior correlation structure</i>
1	global	/	/	/	α_j	iid
/	state	state_id	l	$\{1, \dots, 54\}$	λ_{sj}	spatial (BYM2)
/	individual	age_id	a	$\{1, \dots, 6\}$	η_{aj}^A	random-walk
/		income_id	h	$\{1, \dots, 5\}$	η_{hj}^H	random-walk
/		sex_id	g	$\{1, 2\}$	γ_{gj}^G	unstructured + shared variance
/		race_id	r	$\{1, \dots, 6\}$	γ_{rj}^R	unstructured + shared variance
/		vote20_id	v	$\{1, \dots, 5\}$	γ_{vj}^V	unstructured + shared variance
z_1	state	2020 R share	/	\mathbb{R}	$\beta_{1j=R}$	iid
z_2		On ballot: R.F.K. Jr.			$\beta_{1j=K}$	
z_3		On ballot: Jill Stein			$\beta_{1j=G}$	
z_4		2020 G share			$\beta_{2j=G}$	
z_5		On ballot: Chase Oliver			$\beta_{1j=L}$	
z_6		2020 L share			$\beta_{2j=L}$	
z_7		On ballot: Cornel West			$\beta_{1j=W}$	
z_8		2020 “stay home” share			$\beta_{1j=\text{stay_home}}$	
/	individual : state	vote20_id \times 2020 j share	/	/	$\zeta_{v,j}$	unstructured + shared variance

Table 2: Predictors and Parameters for the 2024 vote-choice model. ‘iid’ refers to fully independent parameters, or ‘fixed’ effects [22]. ‘unstructured + shared variance’ priors refers to classic random-intercepts. Note: the Democrat choice “D” is taken as the reference category, hence it has no associated predictor.

I regress the dependent variable, which is assigned a categorical likelihood with SoftMax link, onto sex, age, ethnicity, household income and 2020 vote. Sex and ethnicity effects are estimated as unstructured random effects; state¹⁰ effects are assigned a BYM2 prior [16, 43, 8]; income and age effects are given random-walk priors. Separate area-level predictors are created for each dependent variable of interest. Table 2 presents the covariates and parameters used in the model for 2024 vote choice.

¹⁰Because I have an interest in being able to estimate the number of electoral votes won by each candidate, I treat the congressional districts of Nebraska and Maine as separate states.

I present the full Hierarchical Bayesian model below – see [14] for a more attentive explanation of each model component. I describe the generation of given choice $j \in \{1, \dots, J\}$, made by a sampled user $i \in \{1, \dots, n\}$, as follows:

$$\begin{aligned}
y_i &\sim \text{Categorical}(\pi_{i1}, \dots, \pi_{iJ}) && \text{likelihood} \\
\pi_{ij} &= \frac{\exp(\mu_{ij})}{\sum_j \exp(\mu_{i,j})}; && \text{softmax link} \\
\mu_{ij} &= \alpha_j + && \\
&\quad \lambda_{\text{state_id}[i],j} + && \\
&\quad \eta_{\text{age_id}[i],j}^A + \eta_{\text{income_id}[i],j}^H + && \\
&\quad \gamma_{\text{sex_id}[i],j}^G + \gamma_{\text{race_id}[i],j}^R + \gamma_{\text{vote20_id}[i],j}^V + && \\
&\quad \sum_{k_j} \beta_{\{k_j,j\}} \times z_{\{\text{state_id}[i],k_j\}} + && \text{state-level predictor} \\
&\quad \zeta_{\{\text{vote20_id}[i],j\}} \times \nu_{\{\text{state_id}[i],j\}}; && \text{ind. by state interactions} \\
\alpha_j &\sim N(0, 1); && \text{intercept} \\
\lambda_{sj} &= \sigma_j^\lambda \left(\phi_{sj} \sqrt{(1 - \xi_j)} + \psi_{sj} \sqrt{(\xi_j/\epsilon)} \right); && \text{BYM2 effects} \\
\phi_{sj} &\sim N(0, 1); && \text{unstr. state-level effects} \\
\psi_{sj} \mid \psi_{s'j} &\sim N \left(\frac{\sum_{l' \neq l} \psi_{s'l'j}}{\nu_l}, \frac{1}{\sqrt{\nu_l}} \right); && \text{conditional auto-reg. effects} \\
\xi_j &\sim \text{Beta} \left(\frac{1}{2}, \frac{1}{2} \right); && \text{mixing weights} \\
\sigma_j^\lambda &\sim N^+(0, 1); && \text{state-level scale} \\
\gamma_{uj}^U \mid \gamma_{u-1j}^U \dots \gamma_{1j}^U &\sim N(\gamma_{u-1j}^U, \sigma_j^U), &\quad \forall U \in \{A, H\}; && \text{random walk effects} \\
\gamma_{uj}^U &\sim N(0, \sigma_j^U), &\quad \forall U \in \{G, R, V\}; && \text{unstructured effects} \\
\sigma_j^U &\sim N^+(0, 1); && \text{random effect scales} \\
\beta_j &\sim N(0, 1). && \text{fixed state cov. effects} \\
\zeta_j &\sim N(0, \sigma_j^\zeta). && \text{unstr. ind. by state effects} \\
\sigma_j^\zeta &\sim N^+(0, 1); && \text{ind. by state scale}
\end{aligned}$$

6.2.1 Learning from Stateless Users

PoSSUM’s geographic filtering ensures users who are selected for analysis are most likely based in the US. The protocol does however allow for the inclusion in the sample of users whose state of residency within the US – their 2nd level geography – is unknown. Learning from these users can bring to bear evidence pertaining the relationship between individual-level attributes, such as age, gender, education, past-vote, etc., and 2024 voting preferences.

I consider an approach that uses two separate linear predictors, one for the observations missing a state, and one for those observations which are complete. For the latter, the linear predictor is exactly as described above; for the former, the following linear predictor is used:

$$\begin{aligned}\mu_{ij}^{s'} = & \alpha_j + \\ & \Xi_j + \\ & \lambda_{\text{state_id}[i],j} + \\ & \eta_{\text{age_id}[i],j}^A + \eta_{\text{income_id}[i],j}^H + \\ & \gamma_{\text{sex_id}[i],j}^G + \gamma_{\text{race_id}[i],j}^R + \gamma_{\text{vote20_id}[i],j}^V \\ \Xi \sim & N(0, 1).\end{aligned}$$

where Ξ_j is a ‘no-state’ fixed effect, which captures the average difference between the baseline level of support relative to the average state, and the ‘no-state’ pool’s support. Effectively, I am treating the ‘no-state’ label as a separate, independent state, which is not pooled towards the state-level effect average. The remaining individual-level coefficients are still informed by these users’ preferences. Making out-of-sample predictions I then use the linear predictor for the observations with known states, effectively discarding the no-state effect.

6.2.2 Alleviating Attenuation Bias

Post-stratified state-level estimates of vote share generated from the above model will tend to display attenuation bias, performing poorly in ‘tail areas’. The size of the attenuation bias is roughly proportional to the size of the difference between the true state-level performance of the candidate, and the candidates’ average performance across states.

Attenuation bias arises as a direct result of modeling choices. Parsimonious models typically perform well for the average state, at the cost of large attenuation bias in tail states. Consider as an example the impact of assuming the relationship between individual-level covariates and voting propensity is constant across states. Similarly, regularising coefficients via partial-pooling will foster attenuation bias by design, as tail effects are smoothed towards the size of average effects. As a result, “tail areas”, which are typically the product of various “tail effects”, are themselves smoothed towards the “average state”.

One might therefore be tempted to fit a “deep” model [26], which makes few parametric assumptions and considers the full set of complex interactions to generate state-level estimates. One issue with this approach is computational tractability, which is part of a trilemma: under typical resources constraints, one cannot have all three of a) flexible modeling; b) simulations from well-behaved posterior distributions; c) fast fitting times which make frequent

(e.g. daily) model updates viable. Beyond that, there are concerns related to over-fitting to unrepresentative samples. I have discussed already issues with LLMs' tendency to exaggerate relationships due to reliance on stereotypes, which affects the representative quality of the PoSSUM data, and make structured regularisation desirable. Unobservable selection effects into online samples also play a role.

One way to address the attenuation bias in a parsimonious, conservative, and tractable manner is to selectively relax regularisation using interaction effects between area-level and individual-level covariates. These have to be sensibly informed by prior knowledge in the data generating process. In the context of vote share estimates, we wish to relax the regularisation pressure at the area-level proportionally to the level of attenuation bias, which we expect to be large wherever a given candidate is most / least dominant, relative to their average performance. A good predictor for this expectation is the candidates' past performance in the given area, which is one component of the state-level predictor. We further know that individual-level vote-choice is driven primarily by past-vote effects. By interacting the individual-level past vote and the area-level past vote-share estimate for the candidate at hand, we can achieve the desired effect:

$$\mu_{ij} = \dots + \sum_{k_j} \beta_{\{k_j, j\}} \times z_{\{\text{state_id}[i], k_j\}} + \zeta_{\{\text{vote20_id}[i], j\}} \times \nu_{\{\text{state_id}[i], j\}} + \dots ;$$

where k_j represents the index of a column from the state-level predictor matrix z which is used to predict choice j ; $\beta_{\{k_j, j\}}$ is the fixed effect of that state-level predictor on choice j ; ν is the matrix of past measurements of the dependent variable at the state-level (e.g. 2020 vote share of choice j in a given state), typically a subset of z ; and $\zeta_{\{v, j\}}$ is the gradient of the effect of choice j 's past state-level measure of choosing option j , for an individual who voted for option j in 2020.

The interaction allows us to account for an additional gradient in the effect of individual-level past-vote across states, allowing the state-level post-stratified estimates to escape, to some degree, regression to the mean.

6.2.3 Aggregating Polls

Aggregating PoSSUM samples is possible – this has the potential to improve the accuracy of estimates for a given fieldwork period by leveraging information from previously fielded polls, and discounting these at an exponential rate [38]. This effect is achieved in the above model via the introduction of a fieldwork-date random effect with a random-walk prior:

$$\begin{aligned} \mu_{ij} &= \dots + \eta_{\text{poll_id}[i], j}^P + \dots \\ \eta_{pj}^P \mid \eta_{p-1, j}^P \dots \eta_{1j}^P &\sim N(\eta_{p-1, j}^P, \sigma^P) \\ \sigma^P &\sim N^+(0, 1). \end{aligned}$$

This effect is simplistic, in that it only captures national-level trends across fieldwork days, whilst in reality state trends over a campaign can vary. Larger state-level sample sizes could justify the inclusion of a more comprehensive state-level random walk with an informative covariance prior [30].

6.3 Estimation

The Hierarchical Bayesian model is implemented in **Stan**¹¹ [11, 59]. To encourage well behaved MCMC sampling, covariates are standardised and random effects are estimated via non-centered parametrisation [48]. The BYM2 effects are specified according to non-generative improper priors [44, 16]. General recommendations for weakly-informative priors [20] are adhered to.

Models are fit separately for each of the 5 polls fielded over the course of the election. A model including fieldwork-period random effects is also fit to an aggregate dataset obtained by stacking the silicon samples.

Each model is ran for 8 chains of 5,000 iterations, with the first 4,750 used as burn-in. The chains are thinned by a factor of 4 to account for auto-correlation, and the `max_treedepth` parameter is set to 15 to allow convergence of the otherwise problematic spatial structure. This results in around 500 posterior samples from each model.

Posterior Predictive Distribution of Crosstabs Posterior simulations for each model parameter can be used to generate samples from the posterior predictive distribution of cell-level choice probabilities. Letting c index the rows of the stratification frame, each representing a population cell, the posterior choice probability samples are derived as follows:

$$\begin{aligned}\tilde{\mu}_{scj} = & \tilde{\alpha}_{sj} + \\ & \tilde{\lambda}_{s, \text{state_id}[c], j} + \\ & \tilde{\eta}_{s, \text{age_id}[c], j}^A + \tilde{\eta}_{s, \text{income_id}[c], j}^H + \\ & \tilde{\gamma}_{s, \text{sex_id}[c], j}^G + \tilde{\gamma}_{s, \text{race_id}[c], j}^R + \tilde{\gamma}_{s, \text{vote20_id}[c], j}^V + \\ & \sum_{k_j} \tilde{\beta}_{\{s, k_j, j\}} \times z_{\{\text{state_id}[c], k_j\}} + \\ & \tilde{\zeta}_{\{s, \text{vote20_id}[c], j\}} \times \nu_{\{\text{state_id}[c], j\}};\end{aligned}$$

$$\tilde{\pi}_{scj} = \frac{\exp(\tilde{\mu}_{scj})}{\sum_j \exp(\tilde{\mu}_{scj})}.$$

Post-stratification is then necessary to obtain population-, crosstab- and state-level estimates of choice-shares. Let $\mathcal{F}(c) = f$ be a function assigning each cell c to a crosstab index f – each c lies in exactly one crosstab, and the crosstabs are exhaustive and mutually exclusive. At the population-level, every cell belongs to the same crosstab $F(c) = f^* \forall c$. It follows that posterior samples from a post-stratified choice-share estimate, for a given crosstab f , are obtained as follows:

$$\tilde{\pi}_{sfj} = \frac{\sum_{c: F(c)=f} \tilde{\omega}_{sc} \tilde{\pi}_{scj}}{\sum_{c: F(c)=f} \tilde{\omega}_{sc}};$$

where ω_{sc} represents a sample from the posterior predictive distribution of stratification-frame weights (obtained from the model described in Section 6.1), and π_{scj} is a posterior sample from the cell-level choice probabilities described above.

¹¹The code for every **Stan** model implemented in this paper is available on the **GitHub** repository.

7 Evaluation

Pre-election opinion polls are useful if they provide an accurate snapshot of the distribution of winning odds across candidates. I score PoSSUM’s **electoral predictive power** by comparing estimates against observed returns at the electoral college district¹² level. The evaluation framework involves statistical measures that reflect essential dimensions of PoSSUM’s estimate quality: i. average error at the point-estimate is captured by $\text{Bias} = \frac{1}{n} \sum_i (\hat{y}_i - y_i)$; ii. the average size of the error is given by $\text{RMSE} = \sqrt{\frac{1}{n} \sum_i (y_i - \hat{y}_i)^2}$; iii. the silicon estimates’ ability to correctly order the observations on a line is measured by the Spearman’s rank

correlation coefficient $\rho = \frac{\sum_i (R_{\hat{y}_i} - \bar{R}_{\hat{y}}) (R_{y_i} - \bar{R}_y)}{\sqrt{\sum_i (R_{\hat{y}_i} - \bar{R}_{\hat{y}})^2} \sqrt{\sum_i (R_{y_i} - \bar{R}_y)^2}}$, where R denotes the rank of a given

observation / prediction; iv. $\text{Coverage}_{90\%} = \frac{1}{n} \sum_{i=1}^n [\hat{y}_i^{5\%} \leq y_i \leq \hat{y}_i^{95\%}]$ measures the calibration of the estimates’ credibility intervals. To measure the agreement between PoSSUM’s estimated silicon density and that implied by traditional polls I use the overlap coefficient $\text{OVL}[p, \hat{p}] = \int_{\mathcal{Y}} \min(p(y), \hat{p}(y)) dy$.

Beyond the ability to predict election results, PoSSUM seeks to enable dynamic measurement of public opinion. To establish PoSSUM amongst state-of-the-arts polling methodologies I must prove it enables *novel learning* and that its estimates are *human-aligned* and *time-sensitive*. I will present an assessment of novelty and alignment exclusively for the last pre-election PoSSUM poll, upon which I build my final pre-election predictions, and use the full set of polls fielded during the campaign to establish time sensitivity.

Crosstab-level Comparisons I wish to establish PoSSUM’s properties at multiple levels of analysis (e.g. by age, gender, education, etc.). Exact election returns at the crosstab level are unavailable. To remedy this, during the campaign I collect a benchmark crosstab-level dataset of polls. The collection protocol is as follows: i) I monitor <https://projects.fivethirtyeight.com/polls/president-general/2024/> in search for newly published polls during the campaign; ii) a poll is identified as a candidate for collection if it shares at least 1 field-work date with PoSSUM; iii) the poll is discarded if crosstabs are behind a paywall; iv) the poll is discarded unless crosstabs are presented in an easily accessible and transcribable format (either pdf, excel, html, or similar). Amazon Textract is used to convert these pdfs into csv files, and these are manually inspected to correct discrepancies; v) weighted or unweighted counts from the crosstabs are stored – polls that provide percentages without the ability to recover counts are discarded. After harmonisation, enough polling data is available to compare PoSSUM with reference pollsters at the levels of gender, ethnicity, age and education-level. A comparison is also made at the level of 2020 past vote, though here the available polls are far fewer (see Figures C.14 to C.19). The database of reference polls is available on the GitHub repository.

¹²These are the constituencies which independently allocate electoral college votes to elect the US President. They typically correspond to US States, with the exceptions of Nebraska and Maine for which they match congressional districts.

Novel learning is operationalised as the ability to correctly capture the *direction* and *magnitude* of changes over past election results. I limit the analysis to a single outcome – namely the *Republican – Democrat* margin, $\mathcal{M} = \pi_{j=R} - \pi_{j=D}$. Novelty can be assessed at multiple levels of analysis: at the national and area-level, I compare estimates against results known with certainty via official election returns; at the crosstab level (e.g. by age, gender, education, etc.) the expanded Stratification Frame (see Section 6.1) can be aggregated at the relevant levels of analysis to provide 2020 vote benchmarks.

Learning an optimal *direction* of change in public opinion implies the ability to minimise the *probability of misdirection*: take $d = \text{sign}(\mathcal{M}^{2024} - \mathcal{M}^{2020})$ to be the observed direction of change, and $\tilde{d}_s = \text{sign}(\tilde{\mathcal{M}}_s^{2024} - \mathcal{M}^{2020})$ to be a single draw from the posterior distribution of predicted changes. The misdirection probability can be estimated as $\widehat{\Pr}(d \neq \tilde{d}) = |\frac{1}{S} \sum_s H(\tilde{d}_s)|$, where H is the Heaviside step function.

An optimal learner of change must correctly capture direction as well as the *magnitude* of change. The *bias* of the learner’s change prediction provides us with a comprehensive evaluation metric: take $\Delta = \mathcal{M}^{2024} - \mathcal{M}^{2020}$ to be the observed change in the margin; $\hat{\Delta} = \hat{\mathcal{M}}^{2024} - \mathcal{M}^{2020}$ is the point estimate of the predicted change; then $\text{bias} = \hat{\Delta} - \Delta$.

Human alignment can be demonstrated by verifying that the observed election outcomes remain consistent with the predictions derived from the PoSSUM posterior. In other words, the empirical data do not provide sufficient evidence to reject PoSSUM as a plausible data-generating process. We can cast this as a hypothesis test: treat the PoSSUM posterior as the null hypothesis (our best prior guess at the result), and posit as the alternative hypothesis that the true result deviates from that distribution. We then compute a *p-value* – the probability that, under the null distribution, we observe a more statistically extreme value than the election result – to assess how plausible the observed result is given PoSSUM’s posterior distribution. To handicap the null, I calculate a one-sided *p-value* for the side of the posterior distribution (relative to its median) that contains the observed value. The Monte-Carlo estimate of this p-value is calculated as follows:

$$\hat{p} = \begin{cases} \frac{1}{S} \sum_s [\tilde{\mathcal{M}}_s \geq \mathcal{M}_{\text{obs}}], & \text{if } \mathcal{M}_{\text{obs}} \geq \text{median}\{\mathcal{M}_s, \dots, \mathcal{M}_S\}, \\ \frac{1}{S} \sum_s [\tilde{\mathcal{M}}_s \leq \mathcal{M}_{\text{obs}}], & \text{if } \mathcal{M}_{\text{obs}} < \text{median}\{\mathcal{M}_1, \dots, \mathcal{M}_S\}. \end{cases}$$

Time sensitivity can similarly be established by comparing temporal dynamics in PoSSUM estimates with those measured by the polling average. Again I rely on Spearman’s rank correlation coefficient over time – large and statistically significant correlation with the polling average at multiple levels of analysis would imply the PoSSUM estimates are *time-sensitive*.

Unfortunately, given the setup of this set of PoSSUM polls, I expect the time-sensitivity test to be largely inconclusive at this stage. PoSSUM polls were in the field for just 5 separate instances during the election campaign, a relatively small number to assess temporal congruence with polling-average dynamics. Moreover, the sample size of each PoSSUM poll was relatively small (approximately $n = 1,000$ in each instance) and smoothed over time (fieldwork dates spanning approximately 1 week) making it challenging to establish patterns of over-time change with high confidence. I will therefore present a single measure of temporal congruence

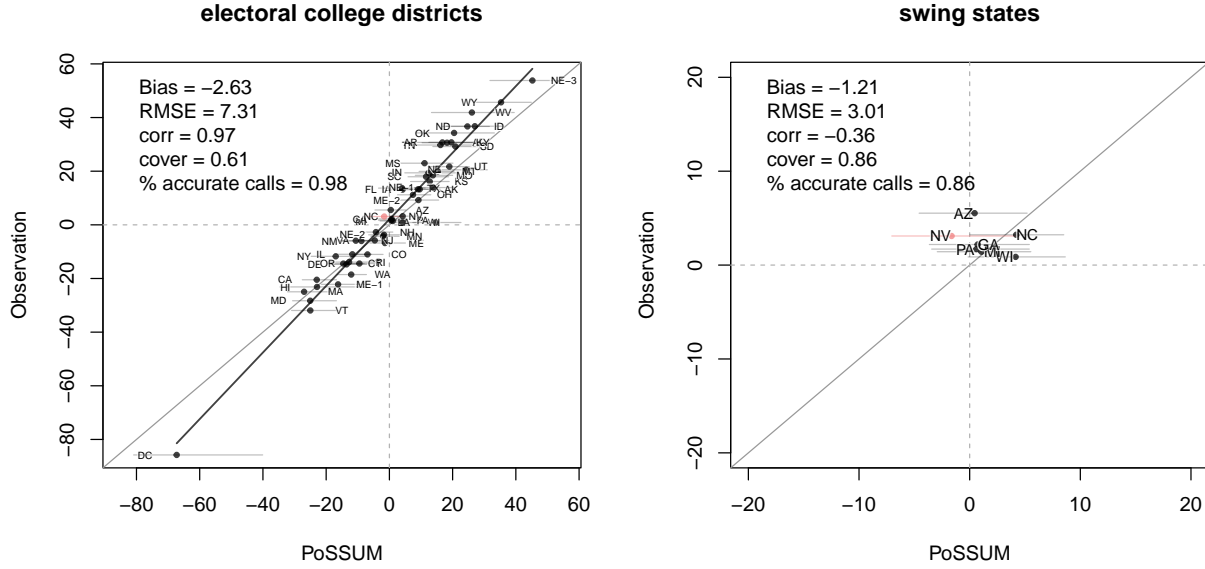


Figure 4: State-level predictive power on the Republican - Democrat margin. Training data includes highly speculative records. Model fit to the final PoSSUM poll, fielded from the 17th to the 26th of October.

between PoSSUM and the polling average at this stage – namely the estimated probability that the temporal correlation between the two is positive – $\widehat{\Pr}(\rho^\tau > 0) = \frac{1}{S} \sum_s [\hat{\rho}_s^\tau > 0]$, where ρ^τ indicates the temporal Spearman rank correlation.

7.1 Electoral Predictive Power

Performance on the $R - D$ Margin Figure 4 presents the comparison of estimates of the *Republican - Democrat* margin from PoSSUM’s last campaign poll, conducted between October 17th and 26th, against the observed state-level election results. The poll is based on a sample of 1,056 synthetic responses – a relatively small sample by traditional standards, and tiny relative to previous MrP efforts on US data [67, 34], it provides above-average informational value.

PoSSUM estimates correctly predicted the Republican win, by calling 98% of electoral college districts correctly – the state of Nevada being the single discordant note. The Spearman rank correlation between point-estimates and observations for the Republican - Democrat margin is generally high, at 0.97, indicating PoSSUM’s ability to accurately order the electoral college districts by the difference in voting intention between Republicans and Democrats. The somewhat elevated RMSE at around 7.3 percentage points emerges as a result of significant attenuation bias, does not affect the ability to call states correctly one way or another, and concentrates the error amongst non-competitive states. On average, margin estimates are characterised by an anti-Trump bias worth 2.6 percentage points, a perhaps surprising finding

given the perception of \mathbb{X} users as being substantially more right-leaning. Coverage of the estimates is unsatisfactorily low, with a mere 61% of observed margins falling in the 90% credibility intervals.

RMSE on the seven competitive swing-states is well contained, at around 3 points – attenuation in less competitive states is responsible for a large overall RMSE. Anti-Trump bias is similarly reduced here to 1.2 points, and is concentrated in the western states of Arizona and Nevada – the crosstab comparisons against state of the arts pollsters in Figure 9 reveal PoSSUM’s inability to capture the Hispanic shift towards Republicans on the last pre-election poll to be the likely cause. Coverage recovers to acceptable levels for these states, with Arizona being the sole state falling outside of the 90% credibility interval. PoSSUM is unable to rank swing states by their Republican - Democrat margin, as evidenced by the negative Spearman correlation coefficient.

Third-Party Under-Performance Figure 5 presents the party-wise vote share estimates, generated from a model fit to an aggregate sample of every poll fielded during the campaign, from August 15th to October 27th, for a total of 4,982 synthetic responses. I use this aggregate sample to attempt a meaningful appraisal of performance on third-parties, whose voters are otherwise too scarce on any given poll to appropriately characterise performance. With the exception of Libertarians, whose state-level vote share is estimated with low RMSE and close-to-perfect coverage, third party prediction poses a challenge for PoSSUM. There is large positive bias in the state-level vote share of RFK Jr. (+3.6%) and Cornel West (+3%). A number of speculative reasons for this over-estimation – such as lack of context available to the LLM with regards to Kennedy’s sui-generis late-endorsement of Trump, or the LLM’s prior beliefs about Dr. West’s chances as a political candidate given his quasi-celebrity status – can be articulated. The systematic partitioning of third-party error according to likely sources is beyond the scope of this paper, though I will discuss these further in Section 8, outlining a clear set of hypotheses to be tested under an experimental setting in future research.

This bias is large enough to have played a consequential role in the predictions. For example in the pivotal state of Wisconsin, where the observed Republican - Democrat margin was a mere 0.8%, and third-parties ended up getting less than 1% of the vote, PoSSUM estimates the margin to be anywhere between 1% and 6.6%. This is a likely consequence of the over-estimation of West and Kennedy, tallying respectively between 1.5% to 4.8%, and between 2% and 5.4%. In this scenario third-parties effectively act to draw votes away from the two major parties. Whilst Kennedy’s vote is generally considered to be contested by both Republicans and Democrats, West’s appeal is limited to Democrats, and as such his over-estimation worked as the equivalent of depressing Democratic turnout.

Comparing with SoTA Pollsters I obtain a collection of state-level polls from the polling aggregator FiveThirtyEight(https://projects.fivethirtyeight.com/polls/data/president_polls_historical.csv). Out of the universe of state-level polls fielded since the 25th of July 2024, I retain the most recent poll related to a Trump v. Harris matchup for each state and pollster. A total of 464 state-level polls remain, fielded by 137 pollsters, for 48 electoral college districts. To perform an apples-to-apples comparison with PoSSUM’s posterior samples for state-level vote-share estimates, and account for their probability distribution, I generate

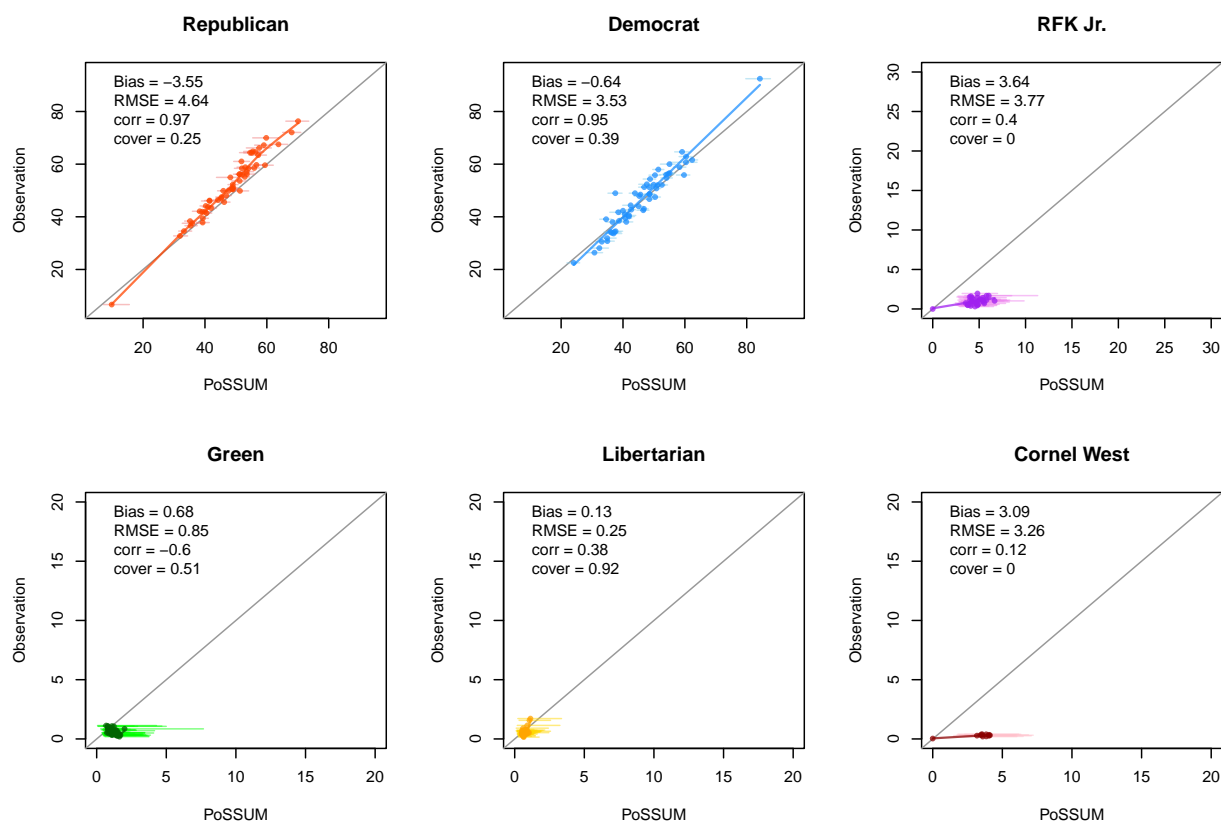


Figure 5: State-level predictive power on vote share by candidate. Training data includes highly speculative records. Model fit to pooled dataset of 5 polls, fielded from the 15th of August to the 26th of October.

an equivalent amount of samples from a Dirichlet distribution (conjugate posterior for the party-choice probability), with $\boldsymbol{\alpha} = N_1 \dots N_J$, where $N_j \forall j \in \{1, \dots, J\}$ is the number of declared supporters of candidate j .

The final comparison involves calculating and comparing the performance metrics described above for the latest pre-election polls fielded by both PoSSUM and the reference pollster. Not every pollster fields in every district, hence the performance comparison will encompass a different set of electoral college districts for each pollster. Presumably pollsters select into fielding in specific states, either because they can assure good performance to clients due to special knowledge, or because the states in question are of broad public interest. Hence the comparison should be broadly favourable to these rival methodologies.

Figure 6 presents the performance difference on the Republican - Democrat margin Δ between PoSSUM and the reference pollsters, by FiveThirtyEight pollster rating. Despite non-negligible anti-Republican Bias (Figure 5), PoSSUM is generally more Republican-leaning than other pollsters. This difference has a clear gradient: pollsters rated lowest by FiveThirtyEight had bias indistinguishable from PoSSUM, whilst state-of-the-arts pollsters significantly favoured Democrats on average. The distribution of this difference is detailed in Figure C.9. The lower anti-Republican Bias translated to PoSSUM obtaining lower average error (RMSE) than state-of-the-arts pollsters, though PoSSUM's RMSE performance was generally indistinguishable from that of the average pollster. Figure 7 presents the distribution of PoSSUM's Δ . A few notable comparisons, ordered in terms of the rival pollster's performance, worst to best: the AI pollster's RMSE on the state-level margin was on average 3.3 points lower than CNN/SSRS; 0.5 points lower than Marist and NYT/Siena; 1.1 points lower than YouGov; 1.8 points lower than Washington Post / George Mason. It was however 0.6 points larger than Morning Consult's; 2.1 points larger than ActiVote; 1.4 points greater than Emerson; 0.7 points greater than Trafalgar; 2.2 points greater than AtlasIntel and 1.4 points greater than Fabrizio/McLaughlin .

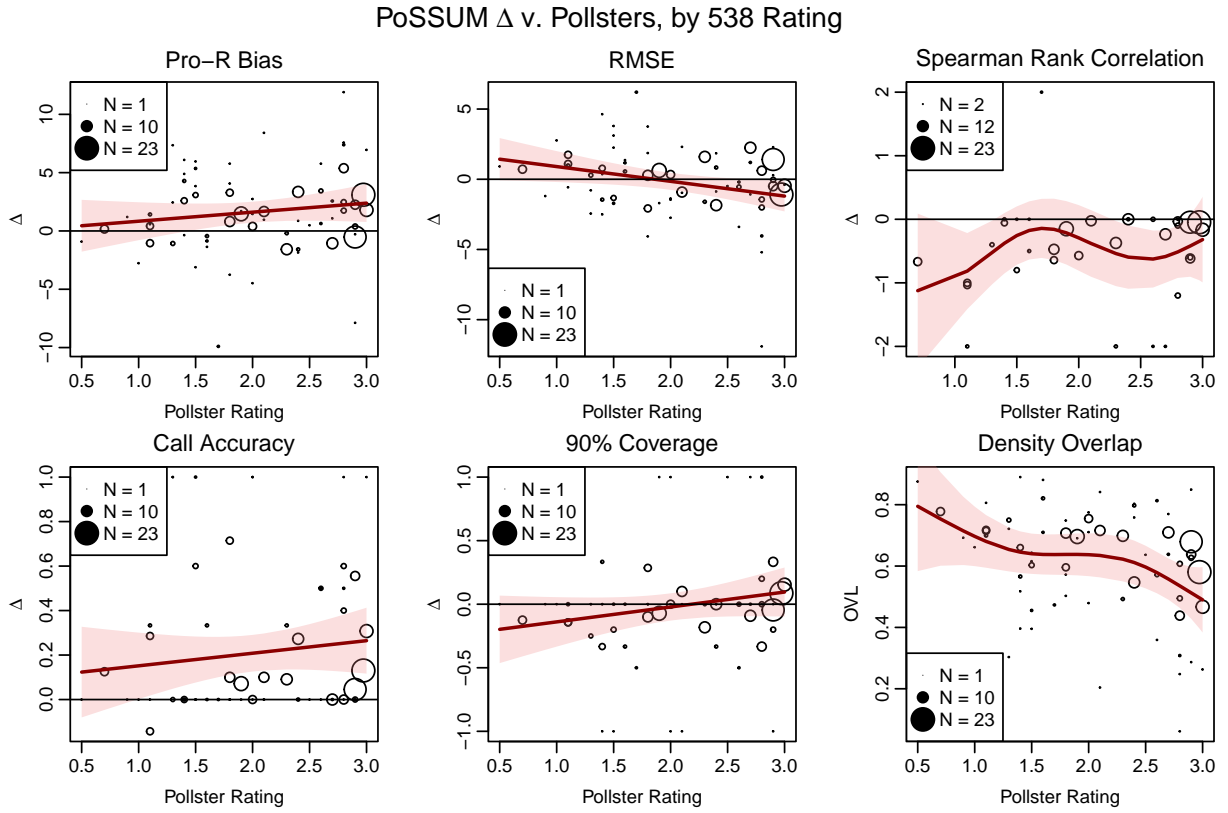


Figure 6: PoSSUM performance difference over reference pollsters ($\Delta = \text{PoSSUM} - \text{Pollster}$) across multiple evaluation metrics by **FiveThirtyEight** pollster ratings. Each comparison is limited to the latest pre-election polls fielded by both parties, in the electoral college districts where the reference pollster has fielded during the campaign.

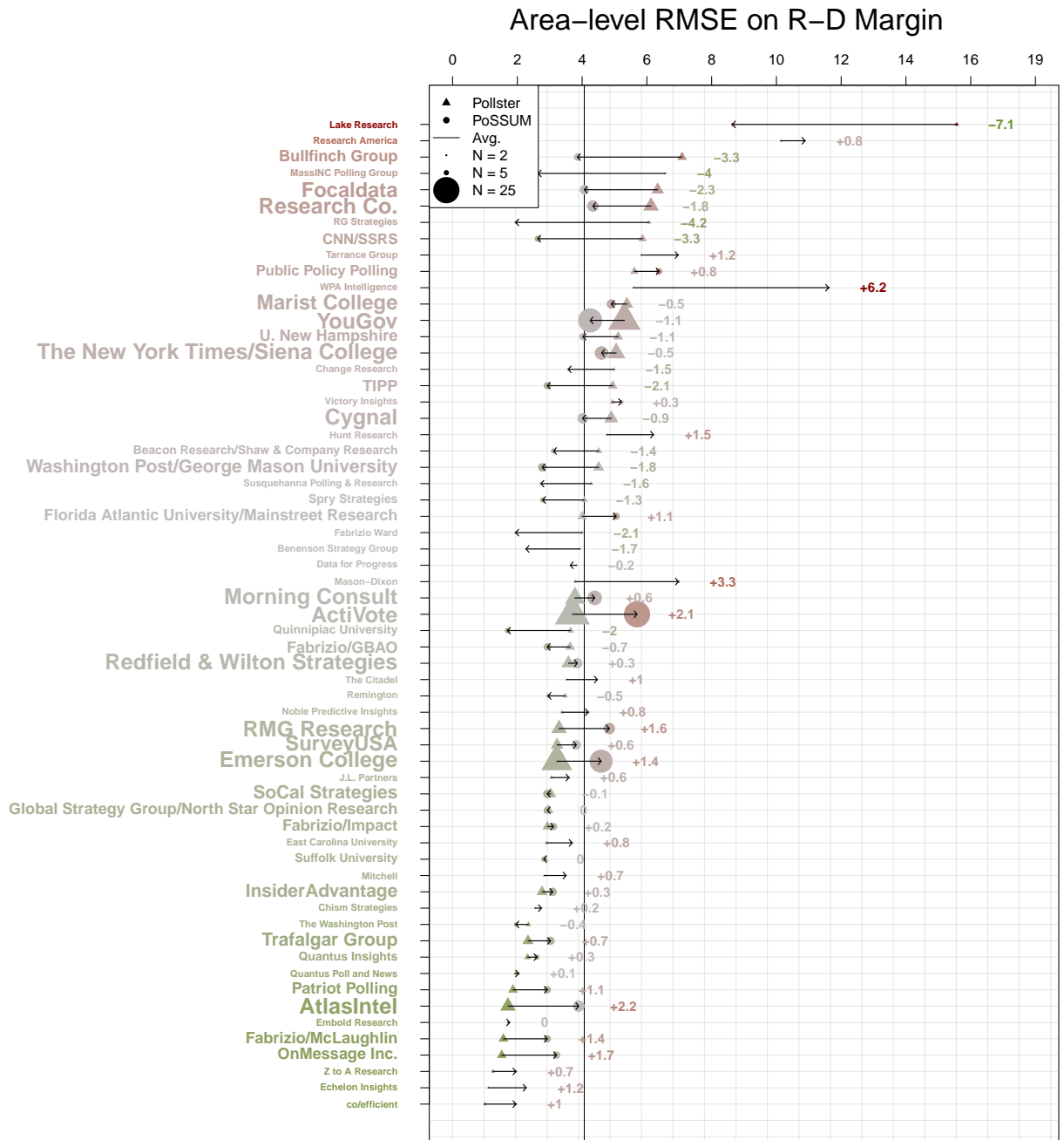


Figure 7: Comparison of RMSE on the area-level Republican margin for each reference pollster (▲) v. PoSSUM (●). Arrow length reflects the difference between estimates, and arrowheads point toward PoSSUM. Pollsters are listed from highest (top) to lowest (bottom) RMSE. The red–green color scale indicates worse–better performance relative to the average. PoSSUM’s RMSE difference (Δ) is displayed to the right of each comparison: green if PoSSUM’s RMSE is lower than the reference pollster, and red if higher. Symbol and label sizes are proportional to the number of areas compared. Only pollsters with data from more than one area are included.

7.2 Novel Learning, Human Alignment and Time-Sensitivity

Figures 8, 9 and Figures C.20 to C.22 in the Appendix present an assessment of the novel learning, human alignment and time-sensitivity of PoSSUM estimates, with respect to observed election results and polling averages. Table 3 summarises these results, according to the metrics presented at the beginning of this section. To interpret the results in context, the table also presents a calculation of the same metrics for state-of-the-arts pollster YouGov.

I find **strong evidence in favour of PoSSUM’s ability to generate novel learning**. The evidence in favour of this is as follows: a) PoSSUM’s point estimate for the change in support is in the correct direction for all crosstabs under considerations, with the exception of those aged 65+ – and even here PoSSUM correctly predicts this is the only crosstab not to experience positive change on the Republican margin; b) the probability of incorrectly predicting the direction of the change for any given crosstab is generally small, proportional to the size of the observed change, and well below random chance for most crosstabs – again with the exception of 65+ voters where we see a toss-up; d) the crosstab-level bias of the predictions is below 4% for most crosstabs, a state-of-the-arts level of performance at the crosstab level. Exceptions to this level of error, varying in degrees of severity, are seen in Hispanics, those in the 18 – 25 and 35 – 44 age brackets, and those who have obtained college degrees.

When comparing PoSSUM and YouGov’s performance on these crosstabs, we see point-estimate error of similar or greater magnitude for YouGov, coupled with a higher frequency of misdirection (largely driven by greater anti-Republican bias in the topline). The higher frequency and sample-size of YouGov polling over the campaign translate to narrower prediction intervals, which leading to more extreme misdirection probabilities – tending to 0 when the point-estimate is in the right direction, and tending to 1 vice-versa. It is notable that YouGov fails to capture any significant dynamics in the 65+ demographic.

There is **moderate evidence favouring human alignment** of PoSSUM’s estimates with observed outcomes. At the 10% significance level, we fail to reject the PoSSUM posterior for all crosstabs except Hispanics, the 18 – 25 and 25 – 44 age brackets, and college-educated voters; the national topline estimate similarly appears statistically different at the 10% level. By contrast, at the more relaxed 5% threshold, we find that both the national topline and the 35 – 44 bracket remain statistically compatible with the PoSSUM distribution.

In absolute terms, most of the PoSSUM crosstab estimates are closely aligned with the observed data. Notable exceptions include the topline result and the 35 – 44 cohort, which lie near the extreme edge of plausibility, and the observed shifts in Hispanic and 18 – 25 voters (albeit with caveats), which appear misaligned. Comparing PoSSUM to YouGov, the former generally yields larger p -values—only the Female and Hispanic crosstabs produce lower p -values under YouGov—suggesting PoSSUM’s estimates are more congruent with actual results. Indeed, the observed margins for topline, Black, 35 – 44, and 65+ groups fall completely outside YouGov’s plausible range. In a complex electoral context, PoSSUM’s learned dynamics thus appear more closely tethered to human beings on the ground than those derived from YouGov’s methods.

The test aimed at **establishing the time-sensitivity of PoSSUM estimates is incon-**

clusive. At face value, the probability of a positive temporal correlation is below random expectations for nearly every crosstab, except for the Female, Black, and 25 – 34 groups—and of these, the result is decisively positive only for the Black crosstab. By contrast, a similar test applied to the state-of-the-art pollster YouGov shows similarly lackluster performance, with only the Black, Hispanic, 35 – 44, and 65+ estimates exhibiting a positive temporal correlation. Several factors may explain this broadly negative performance. First, the method used to compute the correlation may have been primed to report negative results. Specifically, the reference polling average itself is suspect, inasmuch as it is formed by averaging polls conducted at different points in time, overlapping only partially with PoSSUM’s fieldwork. This procedure likely injects noise and smooths the reference benchmark. Similarly, the benchmark derived from YouGov data is also noisy, arising from aggregating multiple polls during each PoSSUM fieldwork window, with potential dependencies among polls that are assumed to be independent (e.g., YouGov polls for CBS or Yahoo vs. those for The Economist). Moreover, aggregating data based on partial overlap with PoSSUM fieldwork periods likely exacerbates this noise, further smoothing and confounding temporal patterns. Although introducing these noisy estimates does not greatly diminish our capacity to evaluate the final PoSSUM fieldwork snapshot against a polling average or the YouGov measurements, it complicates the assessment of changes over time. The resulting compounded noise in the benchmark introduces uncertainty that impedes a conclusive evaluation of PoSSUM’s time-sensitivity in relation to observed outcomes.

In general, PoSSUM exhibits some difficulties with Hispanic voters, for whom there is a 17.8 percentage-point shortfall in the estimated margin shift, a discrepancy without a clear immediate explanation. The lack of novel learning on Hispanic voters in the final polling wave directly contributed to weaker performance in Nevada and Arizona, whereas more accurately captured swings among Black voters enabled correct predictions in Georgia and North Carolina. By contrast, although the error for the 18 – 25 age bracket is even larger at 19.5 points relative to the polling average, this issue is less concerning. Part of the discrepancy can be attributed to inconsistencies in how various pollsters define and aggregate age categories (see the footnote in Table 3). As illustrated in Figure C.21, many polls combine the 25 – 34 and 18 – 24 brackets, and those that do not show alignment with PoSSUM throughout the campaign. Finally, there are smaller underestimations of the Republican margin for the 35 – 44 bracket and for college-degree voters, at 5.9 and 10 percentage points, respectively. Although these discrepancies lack a clear proximate cause, their magnitudes are less pronounced than those for Hispanics and the 18 – 25 bracket.

<i>Crosstab</i>	Δ_{R-D}^{PoSSUM}	$\Delta_{R-D}^{Avg.}$	$\hat{d} = d$	$\Pr(\hat{d} \neq d)$	<i>bias</i>	<i>p-value</i>	$\Pr(\rho^\tau > 0)$
National	2.5	5.9	✓	0.143	-3.3	0.087	0.179
Male	3.0	5.2	✓	0.127	-2.1	0.238	0.183
Female	2.0	1.3	✓	0.256	0.7	0.393	0.552
White	1.9	3.9	✓	0.222	-2.0	0.214	0.381
Black	8.7	8.4	✓	0.060	0.4	0.462	0.931
Hispanic	1.1	18.9	✓	0.419	-17.8	0.004	0.270
Age 18–25 ¹	1.9	21.4	✓	0.417	-19.5	0.004	0.423
Age 25–34	8.2	11.8	✓	0.026	-3.6	0.183	0.597
Age 35–44	4.4	10.3	✓	0.135	-5.9	0.091	0.435
Age 45–54	2.1	2.3	✓	0.313	-0.2	0.480	0.375
Age 55–64	2.0	0.9	✓	0.292	1.1	0.389	0.440
Age 65+	0.0	-3.0	✗	0.500	3.1	0.175	0.349
No College Degree ²	0.7	0.7	✓	0.403	0.0	0.500	0.464
College Degree	3.5	13.5	✓	0.048	-10.0	0.000	0.226
	Δ_{R-D}^{YouGov}	$\Delta_{R-D}^{Avg.}^3$	$\hat{d} = d$	$\Pr(\hat{d} \neq d)$	<i>bias</i>	<i>p-value</i>	$\Pr(\rho^\tau > 0)$
National	1.9	5.9	✓	0.018	-4.0	0.000	0.377
Male	3.1	5.7	✓	0.067	-2.6	0.119	0.391
Female	1.1	1.4	✓	0.286	-0.3	0.558	0.359
White	2.6	4.3	✓	0.036	-1.7	0.188	0.256
Black	0.5	10.6	✓	0.446	-10.1	0.000	0.746
Hispanic	13.9	20.2	✓	0.000	-6.3	0.069	0.706
Age 18–25 ¹	13.6	23.4	✓	0.000	-9.8	0.006	0.313
Age 35–44	-1.4	13.3	✗	0.692	-14.6	0.000	0.567
Age 45–54	3.8	1.5	✓	0.056	2.3	0.196	0.276
Age 55–64	-0.8	1.4	✗	0.617	-2.3	0.181	0.292
Age 65+	5.2	-5.3	✗	0.966	10.5	0.000	0.569

¹ Note that few pollsters collect data for the 18 – 25 category, and in general age categories are grossly misaligned. The reference polls for PoSSUM’s age categories are picked to include any age category which has midpoint distance which is ≤ 5 years from PoSSUM’s category. For the 18 – 25 category, this produces a severe smoothing bias – if one nets-out polls which also overlap with the 25 – 34 category, the polling average looks much more similar – though still greater on average – than the PoSSUM estimate for this crosstab. YouGov’s natural age categories are [18 – 29, 30 – 44, 45 – 64, 65+], hence here too there is a fair degree of misalignment – no reference comparison for the 25 – 34 category is possible.

² YouGov does not provide a breakdown by college education. No other single alternative pollster is consistent enough in their crosstab reporting to provide a satisfactory amount of comparison points.

³ The polling average against which I compare the YouGov data is calculated excluding YouGov – hence it is slightly different from the average I compare the PoSSUM crosstabs against.

Table 3: Assessment of *novel learning* (Direction, Probability of Mis-Direction $\Pr(\hat{d} \neq d)$, Prediction Bias *err*), *human exchangeability* (*p-value* with H_0 : PoSSUM posterior distribution) and *time-sensitivity* (chance of positive correlation across PoSSUM fieldwork dates $\Pr(\rho^\tau > 0)$) in PoSSUM’s final pre-election poll. In **red** are highlighted instances of inability to correctly learn novel preferences (wrong point estimate direction, higher-than-chance probability of mis-direction $\Pr(\hat{d} \neq d) > 0.5$, or prediction bias ± 4), incompatibility with PoSSUM’s posterior distribution ($p - value \leq 0.1$), and worse-than-chance ability capturing of temporal dynamics ($\Pr(\rho^\tau > 0) < 0.5$).

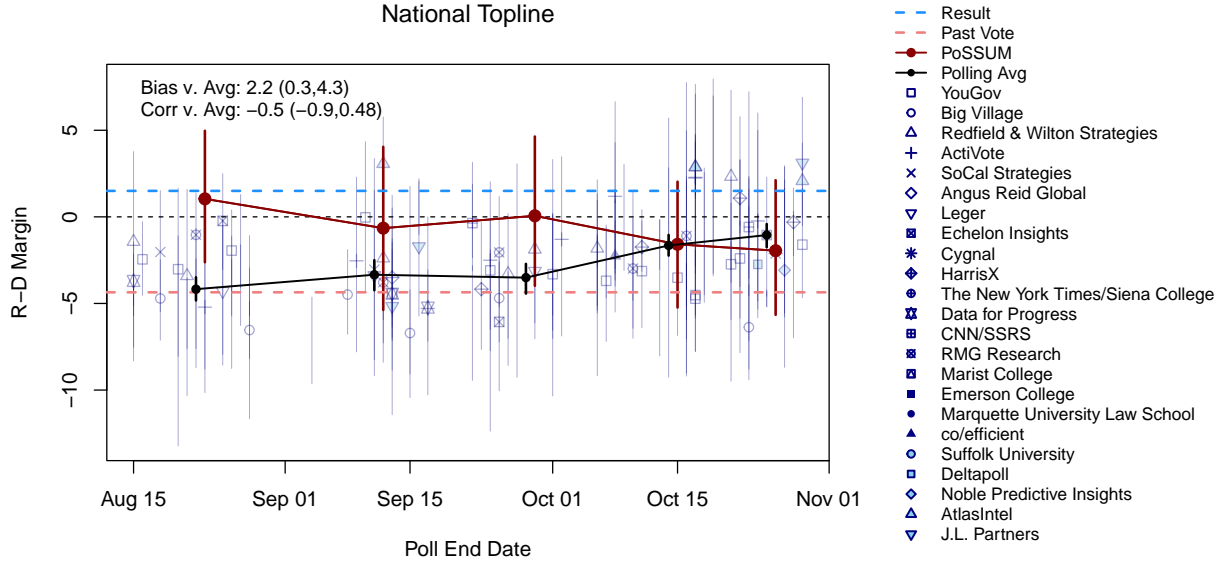


Figure 8: National-level PoSSUM estimates over the course of the campaign, shown alongside individual polls overlapping PoSSUM’s fieldwork periods, the aggregated polling average for each PoSSUM fieldwork window, and the observed 2020 and 2024 outcomes.

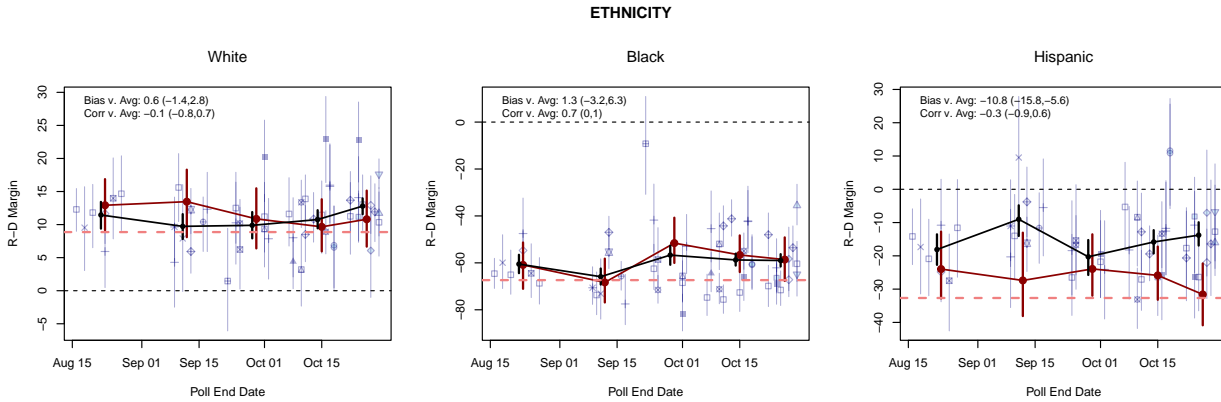


Figure 9: Ethnicity-level PoSSUM estimates over the course of the campaign, shown alongside individual polls overlapping PoSSUM’s fieldwork periods, the aggregated polling average for each PoSSUM fieldwork window, and the reference preferences from 2020.

8 Discussion

This paper has introduced PoSSUM, an end-to-end protocol for unobtrusive polling of social-media users leveraging multimodal LLMs. Tested during the 2024 Presidential Election campaign, it showed that relatively small, calibrated silicon samples — derived from unstructured user data on \mathbb{X} — can produce valid insights into public opinion. PoSSUM accurately predicted state-level outcomes and identified subgroup preferences consistent with official election results and traditional public opinion polls. These findings demonstrate how combining silicon samples with targeted prompting and structured statistical modeling can capture nuanced shifts in public sentiment beyond an LLM’s training cut-off. At the same time, the results highlight persistent challenges, such as underestimating third-party support, reliance on stereotypes and other machine-biases, and the limitations of relying on a single unrepresentative social-media platform. Despite these hurdles, the PoSSUM framework offers a promising path to fully automated, unobtrusive public-opinion polling, capturing granular opinion dynamics of real people without human intervention. The following paragraphs discuss remaining limitations and avenues for future research.

Outstanding Issues Accurately classifying the ethnicity of Hispanic individuals in social media data presents significant challenge. LLMs’ reliance on stereotypical cues [15] may have hindered PoSSUM’s ability to understand the substantial right-shift in this group. It is plausible that right-leaning Hispanics may have been more likely coded as “white”, both due to their lack of stereotypical presentation, and due to the LLMs’ tendency to assume minorities do not support Republicans. Developing more sophisticated methods for identifying the ethnicity of social-media users beyond superficial markers remains an important area of future research.

Regarding educational attainment, a decision was made at the outset of the campaign to substitute college education with income as the optimal socio-economic indicator in both quotas and modeling. This decision was motivated by anecdotal evidence during the piloting phase, which seemed to suggest the LLM overstated the correlation between not holding a college degree and supporting Donald Trump. In hindsight, this choice may have negatively affected model performance, and it is regrettable that a more rigorous comparative assessment was not conducted. Future studies should directly evaluate the LLM biases associated with imputing income v. education levels from social-media data.

The suboptimal performance observed at the extremes of the age distribution appears attributable, at least in part, to limited sampling coverage (Figure A.1). Specifically, PoSSUM exhibited difficulties in adequately populating the youngest and oldest age brackets. One strategy to mitigate this shortfall is to supplement samples from \mathbb{X} with data from other social media platforms. For instance, incorporating TikTok could improve coverage of younger users, whereas Facebook might yield better representation among older populations.

Sources of Third-Party Error Though rigorous experimental evidence beyond the scope of this paper is necessary to unambiguously tease-out the causes of the large third-party miss, a number of competing hypotheses can be set out:

- i. *selection effects* – PoSSUM’s reliance on users discussing politics or trends on \mathbb{X} increases the chance of over-sampling highly engaged users. Queries explicitly seeking third-party

related discussions are liable to inflate the share of third-party voters in the pool, relative to their objective share in the population. The very nature of \mathbb{X} as a “free speech platform” is likely to encourage alternative political discussion;

ii. *parroting of outdated training data* – the LLM’s own training data likely includes polling data from before and around the training cut-off (October 2023). At that time third-parties looked highly competitive. It follows that the baseline rate of labeling users as likely voters of one candidate or another may be biased towards the polling average at the training cutoff period;

iii. *lack of context*: RFK’s endorsement of Trump was not included in the background prompting module, nor was it to be found in the LLM’s training data. Moreover it was a sui-generis endorsement, in that RFK remained on the ballot in many states. It follows that wherever Kennedy was on the ballot, the LLM took him as a legitimate competitor, whilst in the minds of voters he was not viable;

iv. *failure of contextualising a novel candidate*: for candidates without previous political history or much polling, but with an established track-record of policy positions, we can hypothesise the LLM works by placing the candidate on some latent space to put them in relation to other candidates, and to the kinds of voters who may cast a ballot in their favour. Failure to generalise candidate placement on this latent space can lead to poor estimation of vote-choice preferences.

Systematic testing of the hypotheses above requires an experimental setup to uncover counterfactuals – could the LLM have labeled the users differently under different stimuli ? Back-testing on PoSSUM data would likely be unrepresentative of the LLM’s behaviour during the fieldwork period, owing to the temporal instability of the underlying LLM models [9] – hence we cannot hope to learn much more about PoSSUM’s third-party under-performance using the available data. However the hypotheses outlined above can provide a solid foundation to a study of LLM performance on recognising voters of minor parties in multi-party systems.

LLM Identity & Bias The *neutral annotator* approach taken in this paper is somewhat at odds with the agent-based silicon surveying which dominates the literature. The approach is motivated by seeking out the most objective read of a social media user’s timeline, rather than simulate behaviour anew. Under uncertainty however, synthetic responses for specific users rely on the LLM’s best guess, which is often affected by its underlying personality. It follows that a reasonable critique of the *neutral annotator* approach is that the LLM’s default personality is politically biased [54], and that these biases could affect the annotations. An initial defence against this critique is that neutrality might be the preferable default option under uncertainty, given the absence of similar systematic analyses of bias for alternative LLM personalities (i.e. is the “expert forecaster” personality more or less biased than the default ?). Future work should attempt to optimise prompting architecture to minimise annotation biases correlated with relevant features.

The Pollster’s Critique Traditional public opinion researchers are skeptical of the claim that PoSSUM is meaningfully tethered to real, on-the-ground dynamics. This is despite PoSSUM’s effort to link silicon samples to real-life \mathbb{X} users and the text they generate, and further despite LLMs’ proven track-record in labelling political text [61, 24, 14]. These critics may find it difficult to attribute the model’s accurate prediction of the 2024 election to

mere chance, or to dismiss these findings as statistical / computational artifacts: the intricate patterns of change captured by PoSSUM – in both direction and magnitude of preferences, as well as alignment with observed human attitudes across multiple crosstabs, challenge such skepticism.

A more substantive critique concerns the allocation of “credit” for this success: to what extent did the methodological decisions (i.e., treating \mathbb{X} users as a panel, stratified sampling according to demographic quotas, extracting survey content from each user’s posting history, and applying MrP to the silicon samples) materially influence the protocol’s predictive accuracy? Conversely, it remains possible that the base LLM learnings alone might have produced comparable results, had a simpler demographic-prompting approach (e.g., Argyle et al. [4]) been employed on a somewhat representative sample. This would be troublesome, and certainly worthy of further investigation, as it would undermine the generalisability of the performance observed in this experiment.

A key limitation to the present protocol is its non-repeatability. As a proprietary model, `gpt-4o` has been updated in ways that are not publicly documented, making exact replication of PoSSUM’s 2024 polling exercise effectively impossible. Bisbee [9] has highlighted similar concerns regarding evolving LLM capabilities. Moreover, in any future attempts at replication, the 2024 results may already be incorporated (via hard-coding or other data leakage) into the LLM’s knowledge base, confounding efforts to isolate genuine predictive performance.

To disentangle the various sources of predictive success, future research should more rigorously examine the model’s individual-level ability to reconstruct respondents’ preferences and attributes from digital trace data. One viable approach is to match traditional survey data with digital traces and measure the resultant error rates. Even if some degree of individual-level mismatch persists, this need not be detrimental for every application. For agent-based modeling studies that rely on faithfully simulating individual-level interactions to guarantee realistic emergent phenomena [53], greater fidelity at the respondent level might be crucial. In contrast, when the analytical target is an aggregate or crosstab-level estimate, capturing the broader group dynamics may suffice. Nonetheless, rigorously testing and enhancing alignment at the individual level would allow more granular inference and richer subgroup analyses, thereby expanding the utility of the approach.

References

- [1] J. Achiam, S. Adler, S. Agarwal, L. Ahmad, I. Akkaya, F. L. Aleman, D. Almeida, J. Altenschmidt, S. Altman, S. Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- [2] M. Alizadeh, D. Zare, Z. Samei, M. Alizadeh, M. Kubli, M. Aliahmadi, S. Ebrahimi, and F. Gilardi. Comparing methods for creating a national random sample of twitter users. *Social Network Analysis and Mining*, 14(1):160, 2024.
- [3] L. P. Argyle, C. A. Bail, E. C. Busby, J. R. Gubler, T. Howe, C. Rytting, T. Sorensen, and D. Wingate. Leveraging ai for democratic discourse: Chat interventions can improve online political conversations at scale. *Proceedings of the National Academy of Sciences*, 120(41):e2311627120, 2023.
- [4] L. P. Argyle, E. C. Busby, N. Fulda, J. R. Gubler, C. Rytting, and D. Wingate. Out of one, many: Using language models to simulate human samples. *Political Analysis*, 31(3):337–351, 2023.
- [5] R. Baker, J. M. Brick, N. A. Bates, M. Battaglia, M. P. Couper, J. A. Dever, K. J. Gile, and R. Tourangeau. Summary report of the aapor task force on non-probability sampling. *Journal of survey statistics and methodology*, 1(2):90–143, 2013.
- [6] C. Barrie and J. C.-t. Ho. academictwitter: an r package to access the twitter academic research product track v2 api endpoint. *Journal of Open Source Software*, 6(62):3272, 2021.
- [7] C. Barrie, E. Palaiologou, and P. Törnberg. Prompt stability scoring for text annotation with large language models. *arXiv preprint arXiv:2407.02039*, 2024.
- [8] J. Besag, J. York, and A. Mollié. Bayesian image restoration, with two applications in spatial statistics. *Annals of the institute of statistical mathematics*, 43(1):1–20, 1991.
- [9] J. Bisbee, J. D. Clinton, C. Dorff, B. Kenkel, and J. M. Larson. Synthetic replacements for human survey data? the perils of large language models. *Political Analysis*, pages 1–16, 2023.
- [10] T. B. Brown. Language models are few-shot learners. *arXiv preprint arXiv:2005.14165*, 2020.
- [11] B. Carpenter, A. Gelman, M. D. Hoffman, D. Lee, B. Goodrich, M. Betancourt, M. Brubaker, J. Guo, P. Li, and A. Riddell. Stan: A probabilistic programming language. *Journal of statistical software*, 76(1), 2017.
- [12] R. Cerina. *Measuring and forecasting public opinion with non-representative samples*. PhD thesis, University of Oxford, 2020.
- [13] R. Cerina. Machine learning to augment stratification frames in the context of mrp. In *Proceedings of the ITACOSM 2022 Conference*, Perugia, Italy, 2022. Italian Statistical Society.

- [14] R. Cerina and R. Duch. Artificially intelligent opinion polling. *arXiv preprint arXiv:2309.06029*, 2023.
- [15] R. Choenni, E. Shutova, and R. van Rooij. Stepmothers are mean and academics are pretentious: What do pretrained language models learn about you? *arXiv preprint arXiv:2109.10052*, 2021.
- [16] C. Donegan. Flexible functions for icar, bym, and bym2 models in stan. *GitHub*, 2022. URL <https://github.com/ConnorDonegan/Stan-IAR>.
- [17] C. Doyle. Llms as method actors: A model for prompt engineering and architecture. *arXiv preprint arXiv:2411.05778*, 2024.
- [18] E. Ferrara, O. Varol, C. Davis, F. Menczer, and A. Flammini. The rise of social bots. *Communications of the ACM*, 59(7):96–104, 2016.
- [19] Y. Gao, L. Kennedy, D. Simpson, and A. Gelman. Improving multilevel regression and poststratification with structured priors. *Bayesian Analysis*, 16(3):719, 2021.
- [20] A. Gelman. Prior choice recommendations. *Stan Developer Wiki*, 2024. URL <https://github.com/stan-dev/stan/wiki/Prior-Choice-Recommendations>.
- [21] A. Gelman and T. C. Little. Poststratification into many categories using hierarchical logistic regression. 1997.
- [22] A. Gelman, J. B. Carlin, H. S. Stern, D. B. Dunson, A. Vehtari, and D. B. Rubin. *Bayesian data analysis*. Chapman and Hall/CRC, 2013.
- [23] Y. Ghitza and A. Gelman. Deep interactions with mrp: Election turnout and voting patterns among small electoral subgroups. *American Journal of Political Science*, 57(3):762–776, 2013.
- [24] F. Gilardi, M. Alizadeh, and M. Kubli. Chatgpt outperforms crowd workers for text-annotation tasks. *Proceedings of the National Academy of Sciences*, 120(30):e2305016120, 2023.
- [25] I. Goodfellow, Y. Bengio, and A. Courville. *Deep learning*. MIT press, 2016.
- [26] M. Goplerud. Re-evaluating machine learning for mrp given the comparable performance of (deep) hierarchical models. *American Political Science Review*, 118(1):529–536, 2024.
- [27] J. Gottfried. Americans’ social media use, January 31 2024. URL <https://www.pewresearch.org/internet/2024/01/31/americans-social-media-use/>. Accessed: 2024-08-08.
- [28] A. Graefe, J. S. Armstrong, R. J. Jones Jr, and A. G. Cuzán. Combining forecasts: An application to elections. *International Journal of Forecasting*, 30(1):43–54, 2014.

- [29] A. Graefe, H. Küchenhoff, V. Stierle, and B. Riedl. Limitations of ensemble bayesian model averaging for forecasting social science problems. *International Journal of Forecasting*, 31(3):943–951, 2015.
- [30] M. Heidemanns, A. Gelman, and G. E. Morris. An updated dynamic bayesian forecasting model for the us presidential election. *Harvard Data Science Review*, 2(4):10–1162, 2020.
- [31] M. W. Kearney. rtweet: Collecting and analyzing twitter data. *Journal of open source software*, 4(42):1829, 2019.
- [32] G. King and L. Zeng. Logistic regression in rare events data. *Political analysis*, 9(2):137–163, 2001.
- [33] H. Kwak, C. Lee, H. Park, and S. Moon. What is twitter, a social network or a news media? In *Proceedings of the 19th international conference on World wide web*, pages 591–600, 2010.
- [34] B. E. Lauderdale, D. Bailey, J. Blumenau, and D. Rivers. Model-based pre-election polling for national and sub-national outcomes in the us and uk. *International Journal of Forecasting*, 36(2):399–413, 2020.
- [35] Y. LeCun. Do large language models need sensory grounding for meaning and understanding? spoiler: Yes!, 2023. URL https://drive.google.com/file/d/1BU5bV3X5w65DwSMapKcsr0ZvrMRU_Nbi/view.
- [36] L. Leemann and F. Wasserfallen. Extending the use and prediction precision of subnational public opinion estimation. *American journal of political science*, 61(4):1003–1022, 2017.
- [37] C. Li, J. Wang, Y. Zhang, K. Zhu, W. Hou, J. Lian, F. Luo, Q. Yang, and X. Xie. Large language models understand and can be enhanced by emotional stimuli. *arXiv preprint arXiv:2307.11760*, 2023.
- [38] D. A. Linzer. Dynamic bayesian forecasting of presidential elections in the states. *Journal of the American Statistical Association*, 108(501):124–134, 2013.
- [39] P. Liu, W. Yuan, J. Fu, Z. Jiang, H. Hayashi, and G. Neubig. Pre-train, prompt, and predict: A systematic survey of prompting methods in natural language processing. *ACM Computing Surveys*, 55(9):1–35, 2023.
- [40] M. Malik, H. Lamba, C. Nakos, and J. Pfeffer. Population bias in geotagged tweets. In *proceedings of the international AAAI conference on web and social media*, volume 9, pages 18–27, 2015.
- [41] J. Mellon and C. Prosser. Twitter and facebook are not representative of the general population: Political attitudes and demographics of british social media users. *Research & Politics*, 4(3):2053168017720008, 2017.

- [42] G. E. Morris. Artificial intelligence and "big data" cannot replace public opinion polls. *Elliott's notebook*, 2024. URL <https://gelliottmorris.substack.com/p/artificial-intelligence-and-big-data>.
- [43] M. Morris. Spatial models in stan: Intrinsic auto-regressive models for areal data. *GitHub repository*, 2018.
- [44] M. Morris, K. Wheeler-Martin, D. Simpson, S. J. Mooney, A. Gelman, and C. DiMaggio. Bayesian hierarchical spatial models: Implementing the besag york mollié model in stan. *Spatial and spatio-temporal epidemiology*, 31:100301, 2019.
- [45] A. E. Murr. "wisdom of crowds"? a decentralised election forecasting model that uses citizens' local expectations. *Electoral Studies*, 30(4):771–783, 2011.
- [46] OpenAI. Gpt-4 technical report, 2023.
- [47] OpenAI. Introducing gpt-4o, 2024. URL <https://openai.com/index/hello-gpt-4o/>.
- [48] O. Papaspiliopoulos, G. O. Roberts, and M. Sköld. A general framework for the parametrization of hierarchical models. *Statistical Science*, pages 59–73, 2007.
- [49] D. K. Park, A. Gelman, and J. Bafumi. Bayesian multilevel estimation with poststratification: State-level estimates from national polls. *Political Analysis*, 12(4):375–385, 2004.
- [50] J. S. Park, C. Q. Zou, A. Shaw, B. M. Hill, C. Cai, M. R. Morris, R. Willer, P. Liang, and M. S. Bernstein. Generative agent simulations of 1,000 people. *arXiv preprint arXiv:2411.10109*, 2024.
- [51] J. Pasek and M. J. Pasek. Package 'anesrake'. 2018.
- [52] J. Robinson, C. M. Rytting, and D. Wingate. Leveraging large language models for multiple choice question answering. *arXiv preprint arXiv:2210.12353*, 2022. URL <https://arxiv.org/pdf/2210.12353>.
- [53] G. Rossetti, M. Stella, R. Cazabet, K. Abramski, E. Cau, S. Citraro, A. Failla, R. Improta, V. Morini, and V. Pansanella. Y social: an llm-powered social media digital twin. *arXiv preprint arXiv:2408.00818*, 2024.
- [54] D. Rozado. The political preferences of llms. *arXiv preprint arXiv:2402.01789*, 2024.
- [55] I. Rudnyskiy. *openai: R Wrapper for OpenAI API*, 2023. URL <https://github.com/irudnyts/openai>. R package version 0.4.1, <https://irudnyts.github.io/openai/>.
- [56] N. E. Sanders, A. Ulinich, and B. Schneier. Demonstrations of the potential of ai-based political issue polling. *arXiv preprint arXiv:2307.04781*, 2023.
- [57] B. Schaffner, S. Ansolabehere, and M. Shih. Cooperative Election Study Common Content, 2022, 2023. URL <https://doi.org/10.7910/DVN/PR4L8P>.

- [58] P. Schoenegger, I. Tuminauskaite, P. S. Park, and P. E. Tetlock. Wisdom of the silicon crowd: Llm ensemble prediction capabilities match human crowd accuracy. *arXiv preprint arXiv:2402.19379*, 2024.
- [59] C. Stan Development Team et al. Rstan: the r interface to stan. *R package version*, 2(3), 2018.
- [60] P. Törnberg. Best practices for text annotation with large language models. *arXiv preprint arXiv:2402.05129*, 2024.
- [61] P. Törnberg. Large language models outperform expert coders and supervised classifiers at annotating political social media messages. *Social Science Computer Review*, page 08944393241286471, 2024.
- [62] F. t’Serstevens, R. Cerina, and G. Piccillo. Fake news detection via wisdom of synthetic & representative crowds. *arXiv preprint arXiv:2408.03154*, 2024.
- [63] J. Twyman. Getting it right: Yougov and online survey research in britain. *Journal of Elections, Public Opinion and Parties*, 18(4):343–354, 2008.
- [64] U.S. Census Bureau. American Community Survey, 2021 American Community Survey 5-Year Estimates. U.S. Census Bureau, American Community Survey (ACS), 2021. URL <https://www.census.gov/programs-surveys/acs>. Accessed: 2024-08-27.
- [65] S. Van Buuren. *Flexible imputation of missing data*. CRC press, 2018.
- [66] E. Wallace, K. Xiao, R. Leike, L. Weng, J. Heidecke, and A. Beutel. The instruction hierarchy: Training llms to prioritize privileged instructions. *arXiv preprint arXiv:2404.13208*, 2024.
- [67] W. Wang, D. Rothschild, S. Goel, and A. Gelman. Forecasting elections with non-representative polls. *International Journal of Forecasting*, 31(3):980–991, 2015.
- [68] Y. Wang, S. Zhao, Z. Wang, H. Huang, M. Fan, Y. Zhang, Z. Wang, H. Wang, and T. Liu. Strategic chain-of-thought: Guiding accurate reasoning in llms through strategy elicitation. *arXiv preprint arXiv:2409.03271*, 2024.
- [69] Z. Wang, S. Hale, D. I. Adelani, P. Grabowicz, T. Hartman, F. Flock, and D. Jurgens. Demographic inference and representative population estimates from multilingual social media data. In *The World Wide Web Conference*, pages 2056–2067, 2019.
- [70] J. Wei, Y. Tay, R. Bommasani, C. Raffel, B. Zoph, S. Borgeaud, D. Yogatama, M. Bosma, D. Zhou, D. Metzler, et al. Emergent abilities of large language models. *arXiv preprint arXiv:2206.07682*, 2022.
- [71] J. Wei, X. Wang, D. Schuurmans, M. Bosma, F. Xia, E. Chi, Q. V. Le, D. Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35:24824–24837, 2022.

Appendix

A PoSSUM Routines

Input:

- q : optimal API search queries
- w : weight of each query

Output:

- Υ : user-data object composed of profile info v and tweets \mathcal{T}

Routine *get_pool*:

```
 $K \leftarrow \text{length}(q)$  ; # Get: number of queries
 $\Upsilon \leftarrow \emptyset$  ; # Initialize: empty users object
for  $k = 1$  to  $K$  do
   $(\mathcal{T}, v)_{kt} \leftarrow \mathbb{X}(q_k, w_k)$  ; # Call: sample of tweet-user pairs
   $\Upsilon \leftarrow \Upsilon \cup (\mathcal{T}, v)_{kt}$  ; # Store: newly observed users
end
```

Algorithm A.1: Pseudo-code for the `get_pool` routine.

Listing A.1: Search terms for tweets related to candidates involved in the US 2024 presidential election.

```
1 query <-
2 "("
3   Kamala OR VP OR KamalaHarris OR # Democratic
4   candidate terms
5   MAGA OR Trump ORrealDonaldTrump OR # Republican
6   candidate terms
7   Robert Kennedy OR RFK OR RobertKennedyJr OR RFKJr OR KennedyShanahan24 OR Kennedy24 OR # RFK terms
8   Cornel West OR Dr. West OR CornelWest OR # Cornel West
9   terms
10  Jill Stein OR DrJillStein OR # Green
11  candidate terms
12  ChaseForLiberty # Libertarian
13  candidate terms
14 ")"
15 -from:VP -from:KamalaHarris -from:realDonaldTrump -from:RobertKennedyJr -from:CornelWest
16 -from:DrJillStein -from:ChaseForLiberty # Don't sample
17 candidate profiles
18 "
```

Input:

- τ : temporal filter function
- Υ : users object database
- \mathcal{P}^E : entity filter
- E : list of acceptable entities
- \mathcal{P}^G : geographic filter
- G : list of acceptable geographies
- $\mathcal{F} \leftarrow (\mathcal{F}^x, \mathcal{F}^y)$: list of independent and dependent features
- \mathcal{P}^ϕ : feature extraction prompt
- $(\mathbf{X}^\mathcal{Q}, \omega^*, \omega' = \mathbf{0})$: acceptable features, expected frequency, and sample counter
- m : number of tweets per user

Output:

- \mathbf{Z} : survey object with extracted features

Routine *poll_users*:

```

 $\Upsilon^* \leftarrow \Upsilon[\tau(t) = \text{TRUE} \vee \mathcal{G} \neq \emptyset]$  ;           # Filter: recent + valid location
 $N \leftarrow \text{length}(\Upsilon^*)$  ;                               # Get: number of valid users
 $\mathbf{Z} \leftarrow \emptyset$  ;                                   # Initialize: empty survey object

for  $i = 1$  to  $N$  do
     $e_i \leftarrow \text{GPT}\{\mathcal{P}^E(\Upsilon_i)\}$  ;                 # Call: GPT entity filter
    if  $e_i \in E$  then
         $g_i \leftarrow \text{GPT}\{\mathcal{P}^G(\Upsilon_i)\}$  ;           # Call: GPT geographic filter
        if  $g_i \in G$  then
             $\mathbf{X}_i \leftarrow \text{GPT}\{\mathcal{P}^\phi(\Upsilon_i, \mathcal{F}^x)\}$  ; # Call: GPT quota filter
            if  $\mathbf{X}_i \in \mathbf{X}^\mathcal{Q} \cup \omega'_i < \omega_i^*$  then
                 $\mathcal{T}_i^+ \leftarrow \mathbb{X}(\Upsilon_i)$  ;           # Call: sample last  $m$  tweets
                 $\mathbf{z}_i \leftarrow \text{GPT}\{\mathcal{P}^\phi(\Upsilon_i, \mathcal{T}_i^+, \mathcal{F})\}$  ; # Call: GPT extraction
                 $\mathbf{Z} \leftarrow \mathbf{Z} \cup \mathbf{z}_i$  ;                 # Store: survey object
                 $\omega'_i \leftarrow \omega'_i + 1$  ;               # Update: sample quota counter
            end
        end
    end
end
end

```

Algorithm A.2: Pseudo-code for the *poll_users* routine.

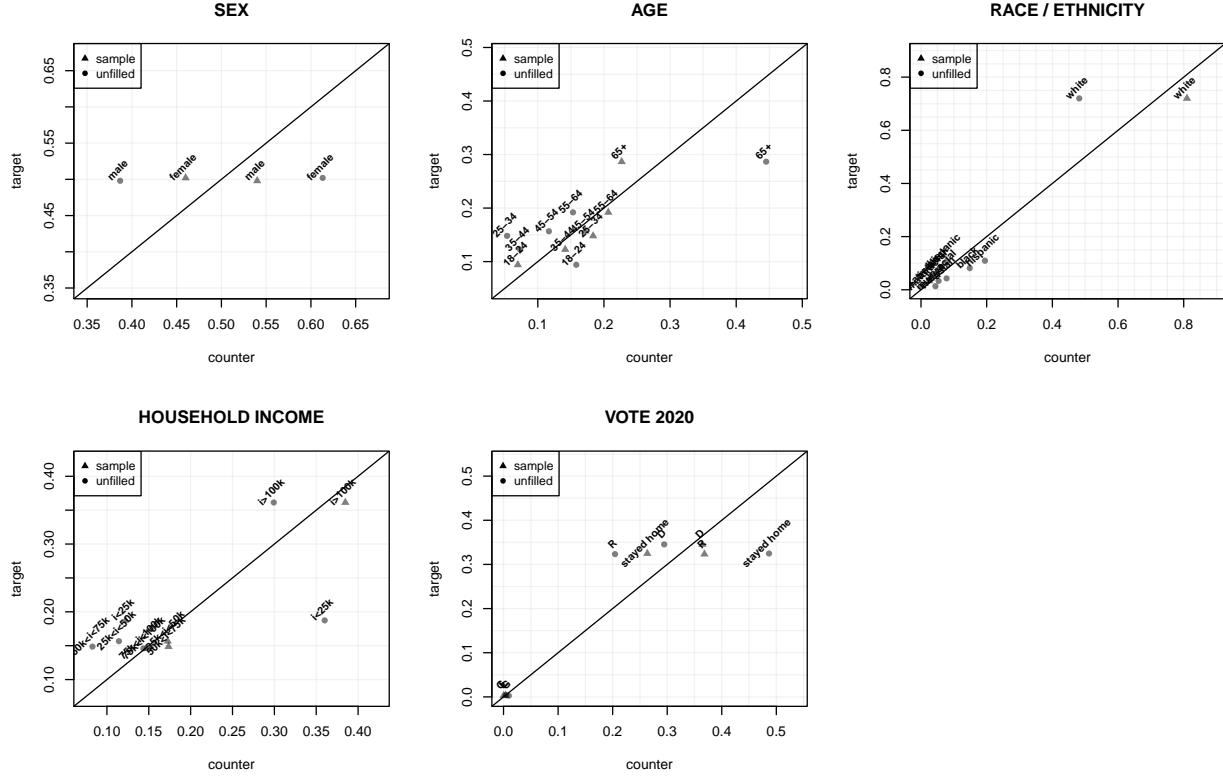


Figure A.1: Sub-group prevalence by filled / unfilled quotas (x-axis) v. target (y-axis), for the final PoSSUM poll, fielded from the 17th to the 26th of October.

Input:

- \mathcal{B} : background information
- \mathcal{F} : features identified for extraction
- \mathcal{M} : mould wrapper
- \mathcal{I} : instructions
- Υ : user data

Output:

- \mathcal{P} : unified prompt

Routine *build_prompt*:

| $\mathcal{P} = \mathcal{B} \parallel \mathcal{M}(\Upsilon) \parallel \mathcal{I}(\mathcal{F})$; # Concatenate: to form unified prompt

Algorithm A.3: Pseudo-code for the *build_prompt* routine.

Input:

- τ : temporal filter function
- Υ : user data
- t : user-wise collection time-stamp

Output:

- Υ^* : filtered user data including only recent users

Routine *temporal_filter*:

```
|  $\Upsilon^* \leftarrow \Upsilon[\tau(t) = \text{TRUE}] ;$            # Filter:  exclude users based on criteria  
    Algorithm A.4: Pseudo-code for the temporal_filter routine.
```

Input:

- Υ : user data
- \mathcal{G} : geographic information associated with users

Output:

- Υ^* : filtered user data with non-null geographic information

Routine *exclude_null_geography*:

```
|  $\Upsilon^* \leftarrow \Upsilon[\mathcal{G} \neq \emptyset] ;$            # Filter:  exclude users with null geography  
    Algorithm A.5: Pseudo-code for the exclude_null_geography routine.
```

Input:

- \mathcal{B} : background information
- Υ : user data
- \mathcal{I}_1 : instructions to generate background-informed features
- \mathcal{F}_1 : empty features object
- \mathcal{I}_2 : operation to extract features
- \mathcal{M} : mould wrapper

Output:

- \mathcal{F}_2 : background-informed features generated by GPT
- \mathbf{X}_i : extracted features for the user data

Routine *background_informed_feature_extraction*:

	$\mathcal{P}_1 = \mathcal{B} \parallel \mathcal{I}_1(\mathcal{F}_1) ;$	# Concatenate: to form 1 st prompt
	$\mathcal{F}_2 \leftarrow \text{GPT}(\mathcal{P}_1) ;$	# Call: GPT to generate bg-informed features
	$\mathcal{P}_2 = \mathcal{B} \parallel \mathcal{M}(\Upsilon) \parallel \mathcal{I}_2(\mathcal{F}_2) ;$	# Concatenate: to form 2 nd prompt
	$\mathbf{X}_i \leftarrow \text{GPT}(\mathcal{P}_2) ;$	# Call: GPT to extract features

Algorithm A.6: Pseudo-code for the background-informed feature extraction routine.

B Prompting Architecture

Listing B.2: Excerpt from the Speculation Module of the prompt, defining how to assign and interpret speculation scores.

```
1 For each selected symbol / category, please note the level of Speculation involved in this selection.
2 Present the Speculation level for each selection on a scale from 0 (not speculative at all, every single element of the
3 user data was useful in the selection) to 100 (fully speculative, there is no information related to this title in the
  user data).
4 Speculation levels should be a direct measure of the amount of useful information available in the user data.
5 Speculation levels pertain only to the information available in the user data -- namely the username, name, description,
  location, profile picture and tweets from this user -- and should not be affected by additional information available
  to you from any other source.
6 To ensure consistency, use the following guidelines to determine speculation levels:
7
8 0-20 (Low speculation): The user data provides clear and direct information relevant to the title. (e.g., explicit mention
  in the profile or tweets)
9 21-40 (Moderate-low speculation): The user data provides indirect but strong indicators relevant to the title. (e.g.,
  context from multiple sources within the profile or tweets)
10 41-60 (Moderate speculation): The user data provides some hints or partial information relevant to the title. (e.g.,
  inferred from user interests or indirect references)
11 61-80 (Moderate-high speculation): The user data provides limited and weak indicators relevant to the title. (e.g., very
  subtle hints or minimal context)
12 81-100 (High speculation): The user data provides no or almost no information relevant to the title. (e.g., assumptions
  based on very general information)
13
14 For each selected category, please explain at length what features of the data contributed to your choice and your
  speculation level.
```

Listing B.3: Example of Multi-feature Object.

```

1 ind.features <- c(
2   'ETHNICITY:
3   E1) white - individuals with origins in any of the original peoples of europe, including, for example, english, german,
      irish, italian, polish, and scottish -- as well as arab or middle-eastern with origins in any of the original peoples
      of the middle east or north africa, including, for example, lebanese, iranian, egyptian, syrian, iraqi, and israeli.
4   E2) black or african american - individuals with origins in any of the black racial groups of africa, including, for
      example, african american, jamaican, haitian, nigerian, ethiopian, and somali.
5   E3) hispanic or latino - includes individuals of mexican, puerto rican, salvadoran, cuban, dominican, guatemalan, and other
      central or south american or spanish culture or origin.
6   E4) asian - individuals with origins in any of the original peoples of central or east asia, southeast asia, or south asia,
      including, for example, chinese, asian indian, filipino, vietnamese, korean, and japanese.
7   E5) american indian or alaskan native or native hawaiian or pacific islander - individuals with origins in any of the
      original peoples of north, central, and south america, including, for example, navajo nation, blackfeet tribe of the
      blackfeet the indian reservation of montana, native village of barrow inupiat traditional government, nome eskimo
      community, aztec, and maya -- as well as individuals with origins in any of the original peoples of hawaii, guam,
      samoa, or other pacific islands, including, for example, native hawaiian, samoan, chamorro, tongan, fijian, and
      marshallese.
8   E6) multiracial - individuals who identify explicitly as belonging to more than one of the racial and ethnic groups above,
      such as biracial individuals with one white and one black parent, or those with a combination of asian and hispanic
      heritage, etc. mixed-race individuals often face unique social experiences, such as celebrating diverse cultural
      holidays, speaking multiple languages, and bridging different cultural perspectives within their families and
      communities.
9   \n',
10  'AGE:
11  A1) under 18 years old
12  A2) 18 to 24 years old
13  A3) 25 to 34 years old
14  A4) 35 to 44 years old
15  A5) 45 to 54 years old
16  A6) 55 to 64 years old
17  A7) 65 or older
18  \n',
19  'SEX:
20  S1) masculine sex - male
21  S2) feminine sex - female
22  \n',
23  'INTEREST IN POLITICS:
24  I1) not interested at all in politics
25  I2) slightly interested in politics
26  I3) moderately interested in politics
27  I4) highly interested in politics
28  \n',
29  'MARITAL STATUS:
30  M1) married - currently legally married and living with a spouse
31  M2) single - never married, including those who are legally separated
32  M3) divorced - legally divorced and not currently remarried
33  M4) widowed - spouse has passed away and not currently remarried
34  \n',
35  "HIGHEST EDUCATIONAL QUALIFICATION:
36  Q1) completed education up to and including high school - high school diploma, vocational training, associate degree
37  Q2) completed education at the college or university level - bachelor's degree, master's degree, doctorate
38  \n",
39  'HOUSEHOLD INCOME BRACKET:
40  H1) up to 25000 USD per year
41  H2) between 25000 and 50000 USD per year
42  H3) between 50000 and 75000 USD per year
43  H4) between 75000 and 100000 USD per year
44  H5) more than 100000 USD per year
45  \n',
46  'GENERAL TRUST IN OTHER PEOPLE:
47  Tru1) always trust other people
48  Tru2) most of the time trust other people
49  Tru3) about half of the time trust other people
50  Tru4) some of the time trust other people
51  Tru5) never trust other people
52  \n',
53  'PAYING ATTENTION TO THE 2024 PRESIDENTIAL ELECTION:
54  Att1) not paying attention at all to the 2024 Presidential election in the US
55  Att2) paying only a little attention to the 2024 Presidential election in the US
56  Att3) paying some attention to the 2024 Presidential election in the US
57  Att4) paying a lot of attention to the 2024 Presidential election in the US
58  \n',
59  .
60  .
61  .
62 )

```

C Supplementary Results

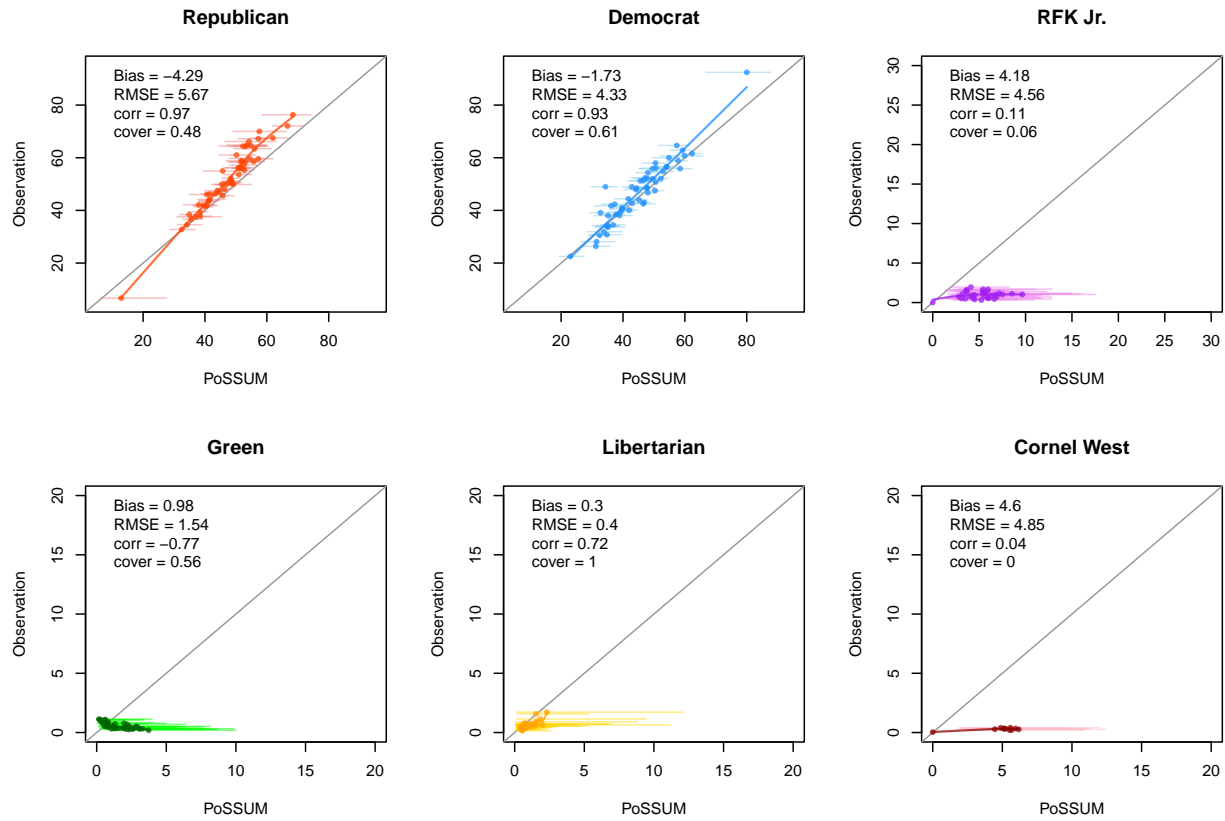


Figure C.1: State-level predictive power on vote share by candidate. Training data includes highly speculative records. Model fit to the final PoSSUM poll, fielded from the 17th to the 26th of October.

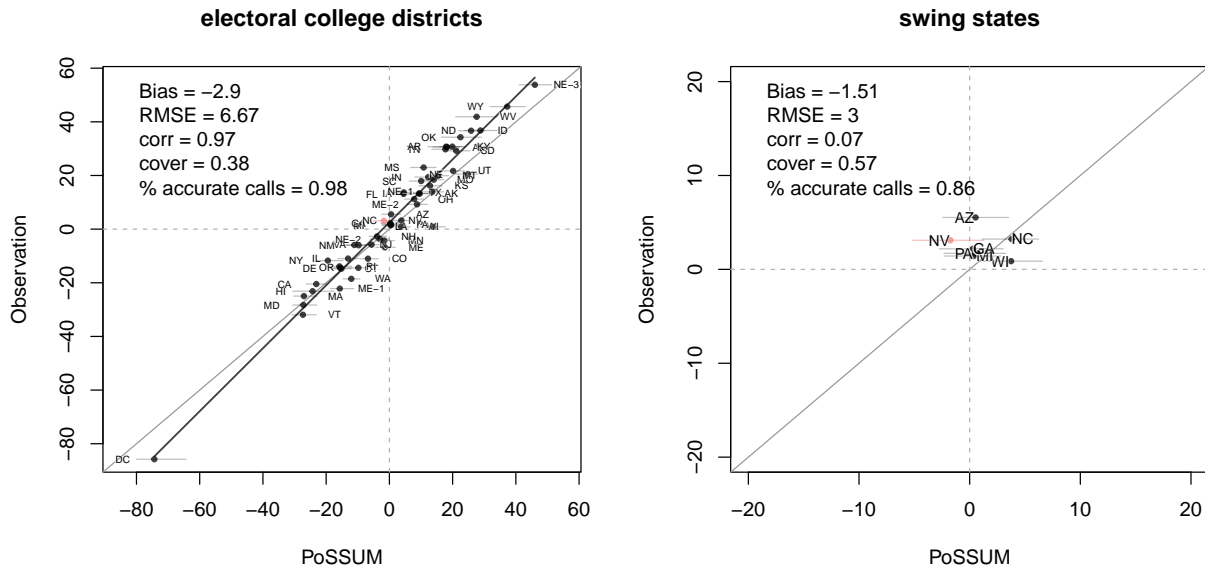


Figure C.2: State-level predictive power on Republican - Democrat margin. Training data includes highly speculative records. Model fit to pooled dataset of 5 polls, fielded from the 15th of August to the 26th of October.

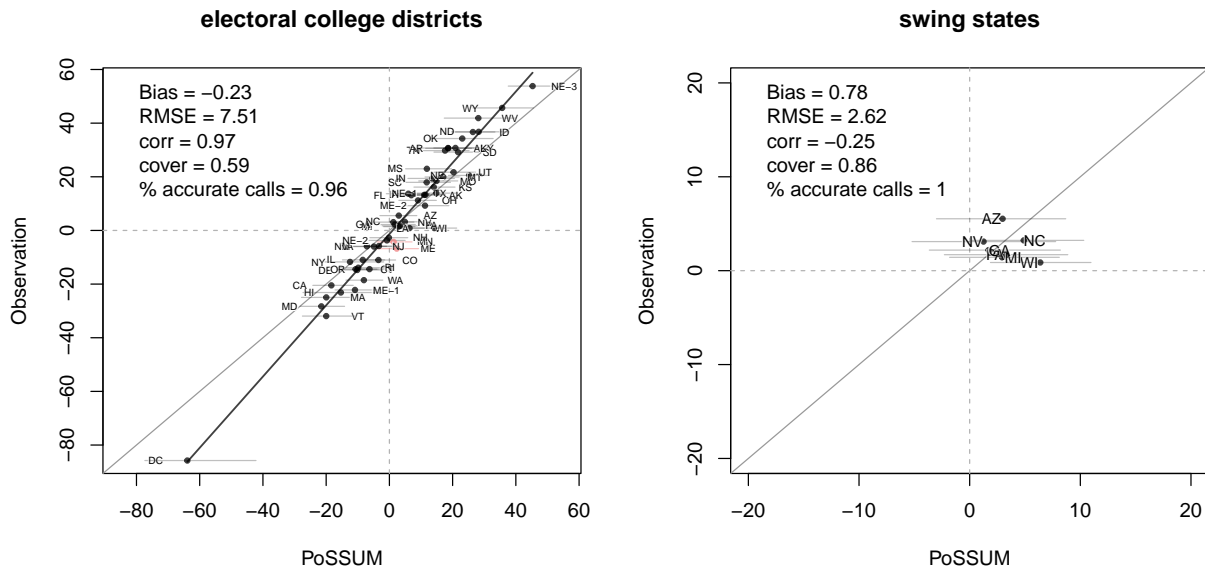


Figure C.3: State-level predictive power on the Republican - Democrat margin. Training data does not include highly speculative records. Model fit to pooled dataset of 5 polls, fielded from the 15th of August to the 26th of October.

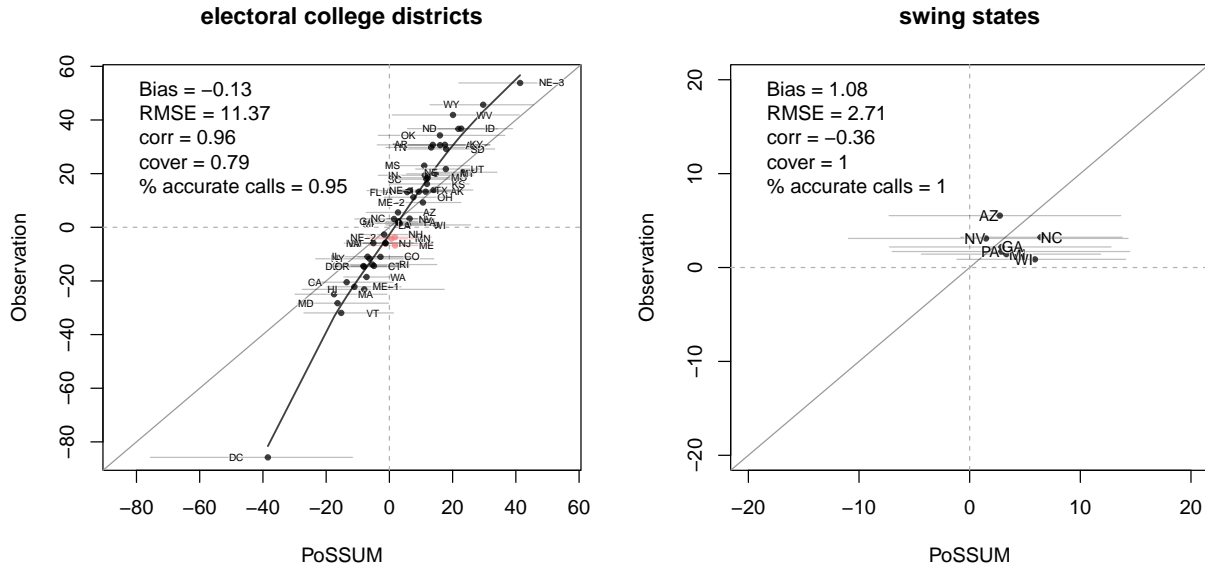


Figure C.4: State-level predictive power on the Republican - Democrat margin. Training data does not include highly speculative records. Model fit to the final PoSSUM poll, fielded from the 17th to the 26th of October.

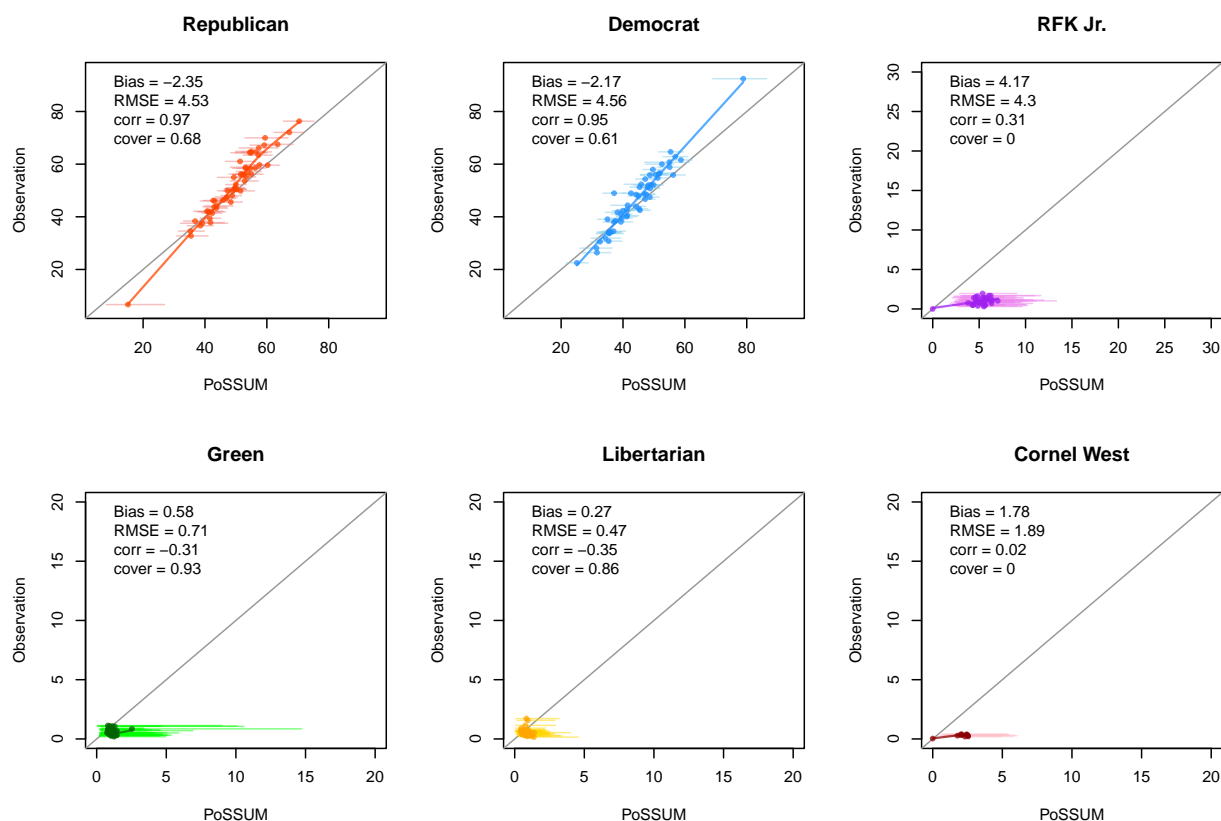


Figure C.5: State-level predictive power on vote share by candidate. Training data does not include highly speculative records. Model fit to pooled dataset of 5 polls, fielded from the 15th of August to the 26th of October.

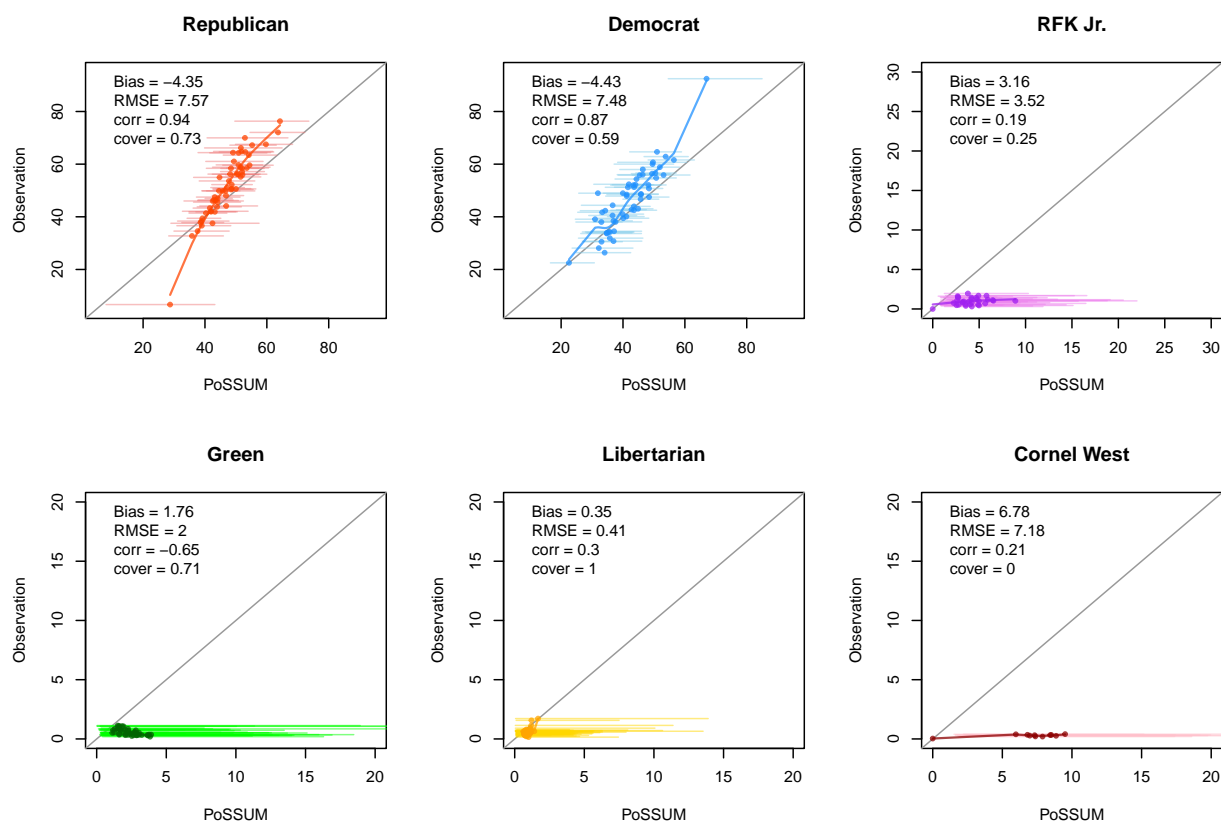


Figure C.6: State-level predictive power on vote share by candidate. Training data does not include highly speculative records. Model fit to the final PoSSUM poll, fielded from the 17th to the 26th of October.

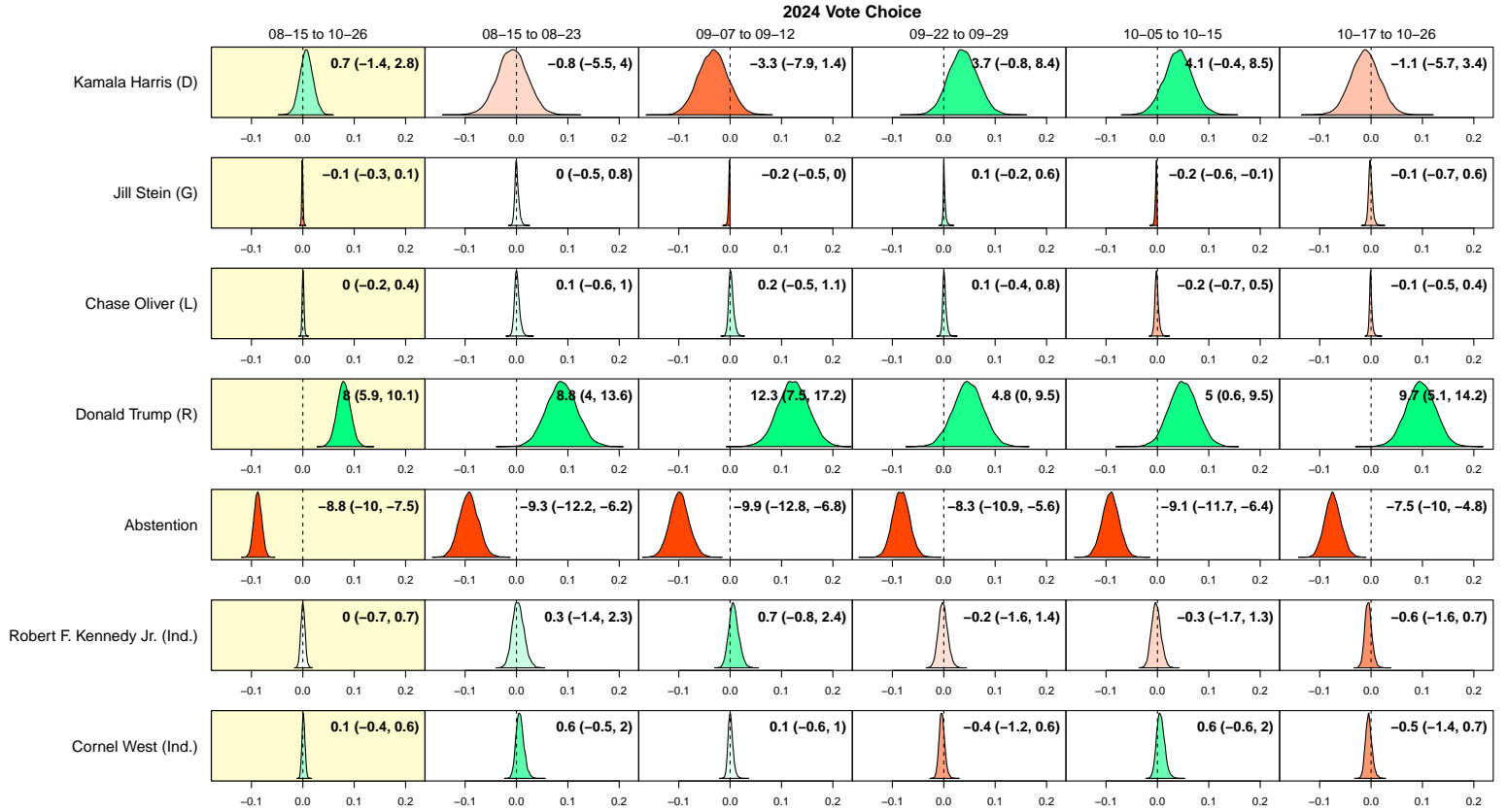


Figure C.7: Posterior distribution of the effect of dropping highly speculative records on the raw training data, for 2024 voting preferences amongst all synthetic samples. The highlighted leftmost column presents the effect on the aggregated complete sample of 5 polls.

Learning from Highly Speculative Records When the MrP model is allowed to learn from highly speculative synthetic records, performance on objective predictive power indicators improves substantially, compared to when this is forbidden. This is demonstrated by a comparison of Figures 4 and C.4.

It is challenging to pin down the exact mechanism through which this works. A substantial decrease in sample size as a result of moderating speculation is surely to be assigned some portion of the blame. Out of the 4,982 users for whom PoSSUM inferred features during the campaign, 2,814 are assigned at least one highly speculative feature necessary to fit the model, and are therefore dropped. The results of this are dramatic levels of attenuation bias (Figure C.4 for a model fit excluding highly speculative records on a single fieldwork period, which enjoys only a tiny state-level sample size of 377 users).

Beyond the sheer effects of sample size, there are also changes in predictive accuracy due to the composition of users which are flagged as highly speculative. Figure C.8 shows the effect of dropping highly speculative users on the distribution of 2024 voting preferences amongst likely voters in the raw data sample. The effect of dropping highly speculative users is to turn the sample of likely voters substantially more Republican. Users with low turnout propensity (Figure C.7), as well as some democrats, are more likely to be the subject of high-levels of speculation. The effect can be seen even after weighting, as shown by lower anti-Republican bias in the MrP estimates trained on the moderately speculative sample. High-levels of speculation are associated with labeling users as medium-to-low income, white ethnicity and abstention in the 2020 election.

The provocative intuition developed from this section is that the LLM’s engagement in speculation is beneficial to preference estimation under the PoSSUM protocol. The fact that speculative records increase sample size does not in itself lead to better estimates – if the estimates were mere noise, we would expect a drop in performance. It follows that the speculative records must contain some relevant auxiliary information, which is not directly acquirable from the underlying mould, but which is useful to address some of the underlying bias in the data. This is evidence in favour of the proposition that LLMs hold information to address bias in unrepresentative samples.

Benefits of Speculation This paper presents a novel prompting strategy to obtain a self-reported *speculation score* from the LLM. Speculation here is defined as the amount of information in the mould which is directly indicative of a category to which the user belongs to. I show that including highly speculative records, defined loosely as records for which at least one of the relevant variables attains a speculation score greater than 80%, is beneficial to the estimation of voting preferences. Future work should focus on uncovering the gradient of this benefit – how sensitive are preferences to varying degrees of speculation ? When the LLM is speculating, it is in effect engaging in a similar process as multiple imputation. Unlike existing imputation algorithms, it is able to bring in knowledge external to the dataset at hand, which it has acquired during the training phase. An interesting question for future research is to identify when this external knowledge outperforms internal knowledge – under what conditions does LLM imputation outperform classic multiple imputation models ?

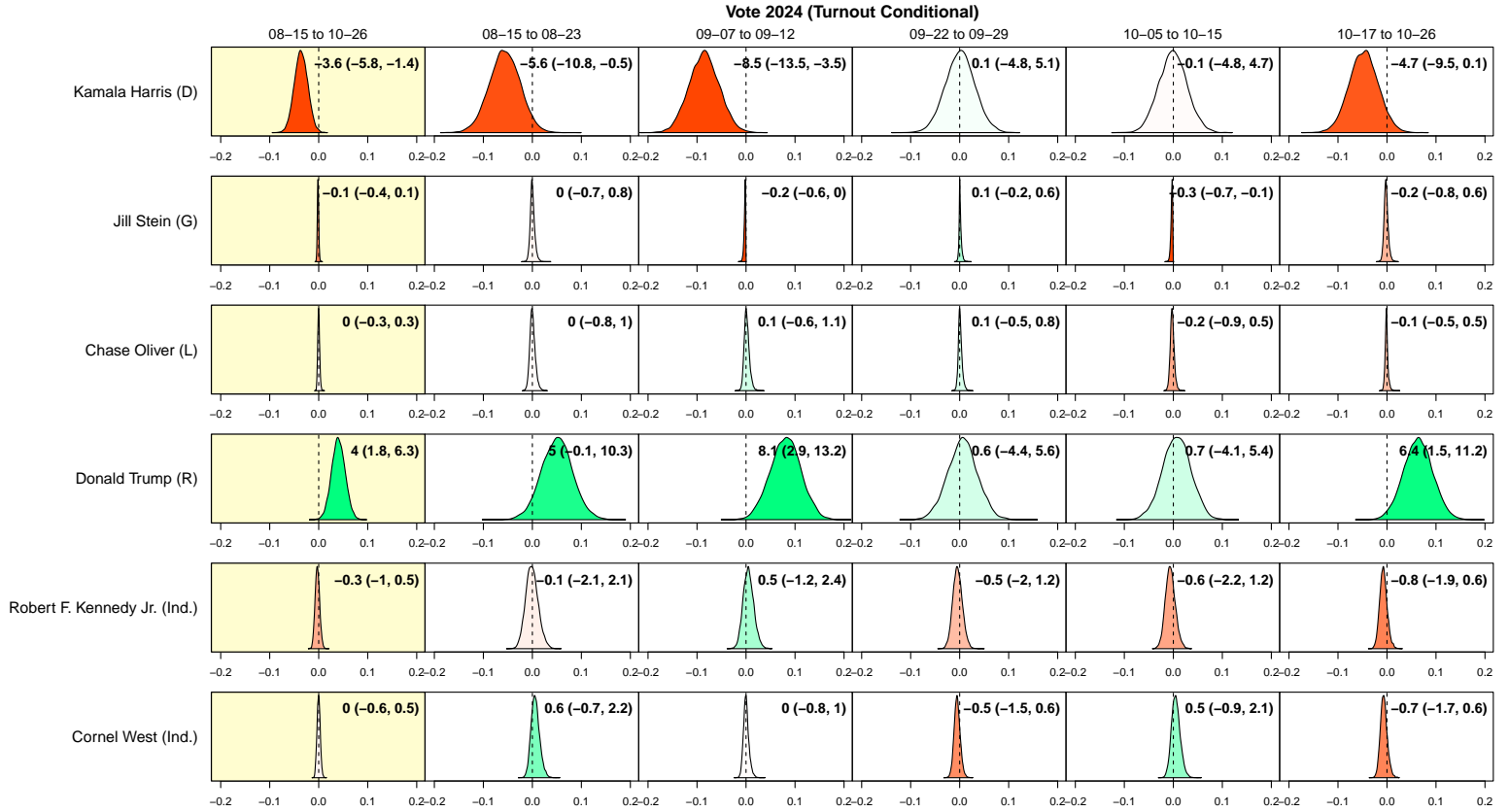


Figure C.8: Posterior distribution of the effect of dropping highly speculative records on the raw training data, for 2024 voting preferences amongst likely voters. The highlighted leftmost column presents the effect on the aggregated complete sample of 5 polls.



Figure C.9: Comparison of Pro-Trump Bias on the area-level Republican margin for each reference pollster (▲) v. PoSSUM (●). Arrow length reflects the difference between estimates, and arrowheads point toward PoSSUM. Pollsters are listed from highest (top) to lowest (bottom) pro-Trump Bias. The blue–red color scale indicates lower–higher pro-Trump Bias relative to the average. PoSSUM’s Bias difference (Δ) is displayed to the right of each comparison: red if PoSSUM favours the Republican more than the reference pollster, blue vice versa. Symbol and label sizes are proportional to the number of areas compared. Only pollsters with data from more than one area are included.



Figure C.10: Comparison of Spearman Rank Correlation on the area-level Republican margin for each reference pollster (▲) v. PoSSUM (●). Arrow length reflects the difference between estimates, and arrowheads point toward PoSSUM. Pollsters are listed from highest (top) to lowest (bottom) Rank Correlation Coefficient. The red–green color scale indicates worse–better performance relative to the average. PoSSUM’s Rank Correlation difference (Δ) is displayed to the right of each comparison: green if PoSSUM’s Rank Correlation is greater than the reference pollster, and red if smaller. Symbol and label sizes are proportional to the number of areas compared. Only pollsters with data from more than 3 areas are included, as Rank Correlation comparisons tend to be unstable below that number.

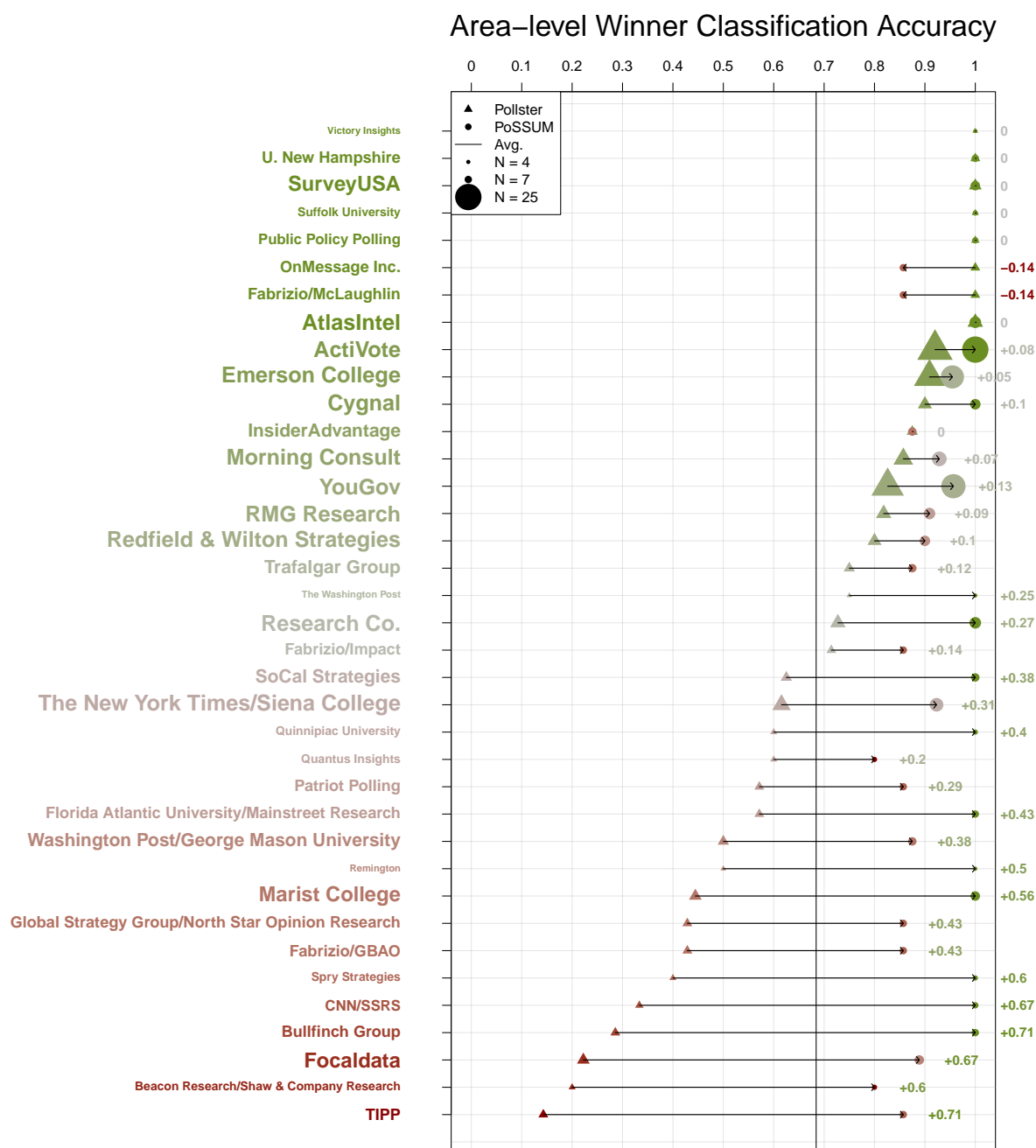


Figure C.11: Comparison of Accuracy on winner prediction at the area-level for each reference pollster (▲) v. PoSSUM (●). Arrow length reflects the difference between estimates, and arrowheads point toward PoSSUM. Pollsters are listed from highest (top) to lowest (bottom) Accuracy. The red–green color scale indicates worse–better performance relative to the average. PoSSUM’s Accuracy difference (Δ) is displayed to the right of each comparison: green if PoSSUM’s Accuracy is greater than the reference pollster’s, and red if lower. Symbol and label sizes are proportional to the number of areas compared. Only pollsters with data from more than 3 areas are included, as Accuracy comparisons to be unstable below that number.



Figure C.12: Comparison of Coverage (90%) on the area-level Republican margin for each reference pollster (▲) v. PoSSUM (●). Arrow length reflects the difference between estimates, and arrowheads point toward PoSSUM. Pollsters are listed from highest (top) to lowest (bottom) Coverage. Stated Coverage is 90%, though over-coverage is not penalised in this comparison. The red–green color scale indicates worse–better performance relative to the average. PoSSUM’s Coverage difference (Δ) is displayed to the right of each comparison: green if PoSSUM’s Coverage is greater than the reference pollster, and red if lower. Symbol and label sizes are proportional to the number of areas compared. Only pollsters with data from more than one area are included.



Figure C.13: Density Overlap Coefficient (OVL) on the area-level Republican margin between each reference pollster and PoSSUM. The red–green color scale indicates lower–higher coverage relative to the average. Symbol and label sizes are proportional to the number of areas compared. Only pollsters with data from more than one area are included.

C.1 Pollster’s Temporal Coverage by Crosstab

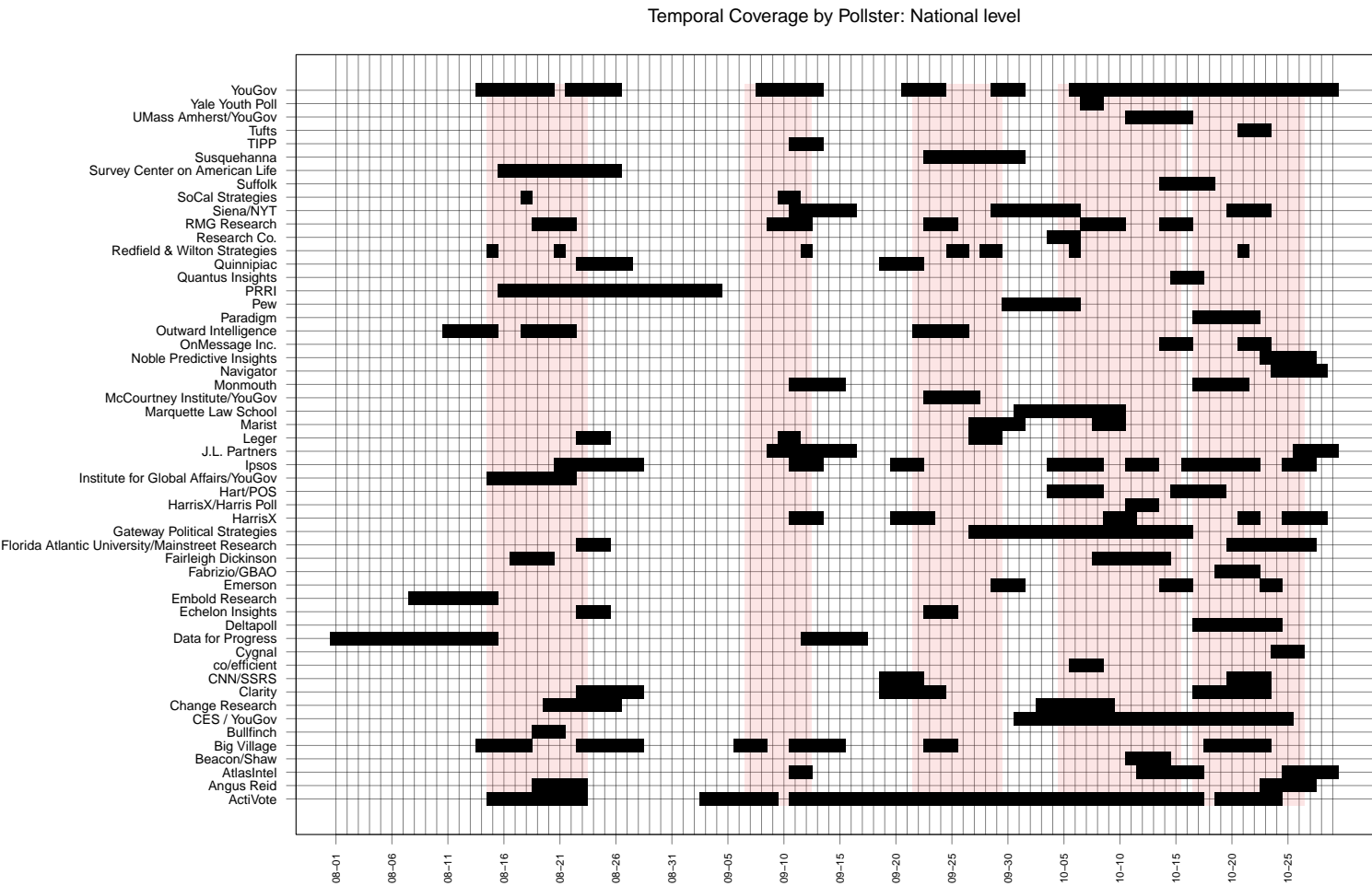


Figure C.14: Temporal coverage of national-level reference pollsters. The light red shading marks the digital fieldwork period for PoSSUM, during which these polls were combined into the average used to evaluate PoSSUM’s results.

Temporal Coverage by Pollster for Crosstab: Gender

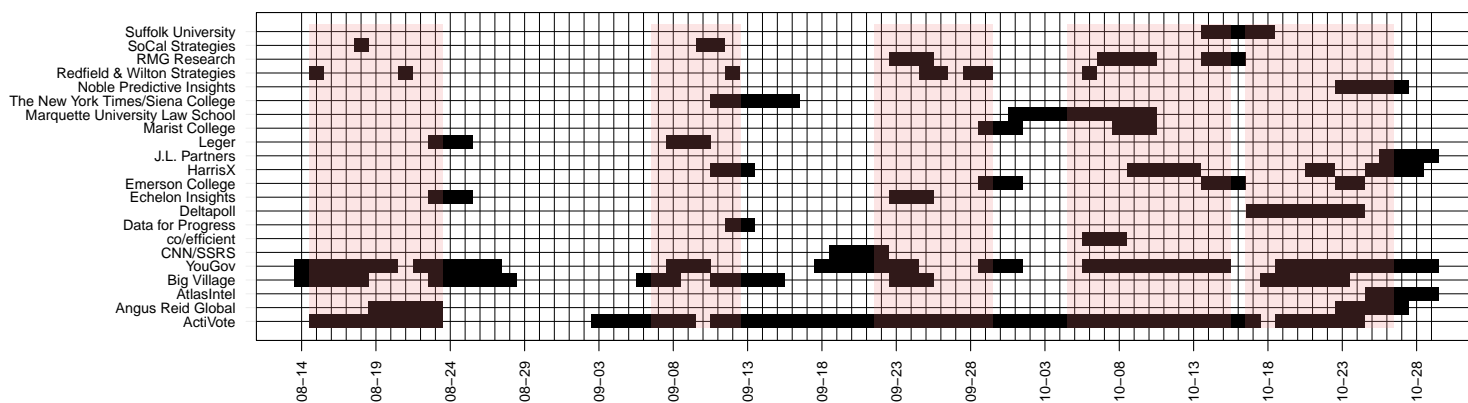


Figure C.15: Temporal coverage of reference pollsters for gender-level crosstabs. The light red shading marks the digital fieldwork period for PoSSUM, during which these polls were combined into the average used to evaluate PoSSUM's results.

Temporal Coverage by Pollster for Crosstab: Race / Ethnicity

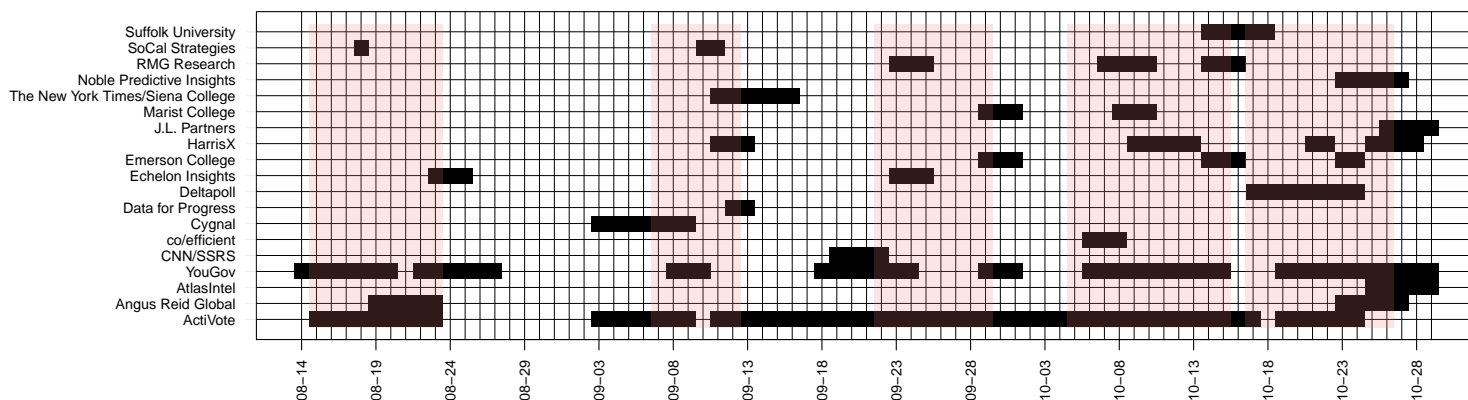


Figure C.16: Temporal coverage of reference pollsters for race / ethnicity -level crosstabs. The light red shading marks the digital fieldwork period for PoSSUM, during which these polls were combined into the average used to evaluate PoSSUM's results.

Temporal Coverage by Pollster for Crosstab: Age

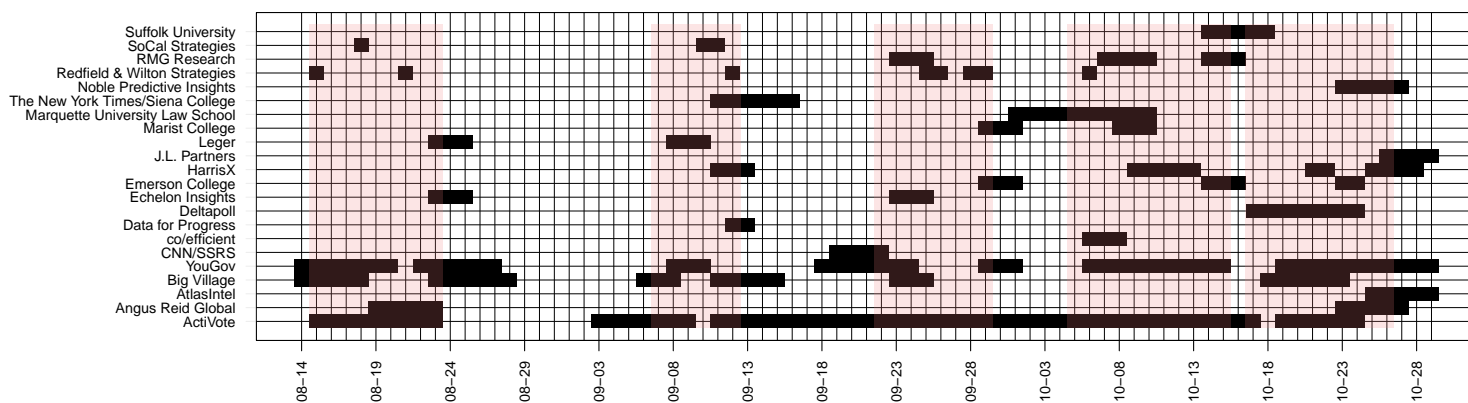


Figure C.17: Temporal coverage of reference pollsters for age-group-level crosstabs. The light red shading marks the digital fieldwork period for PoSSUM, during which these polls were combined into the average used to evaluate PoSSUM's results.

Temporal Coverage by Pollster for Crosstab: Education

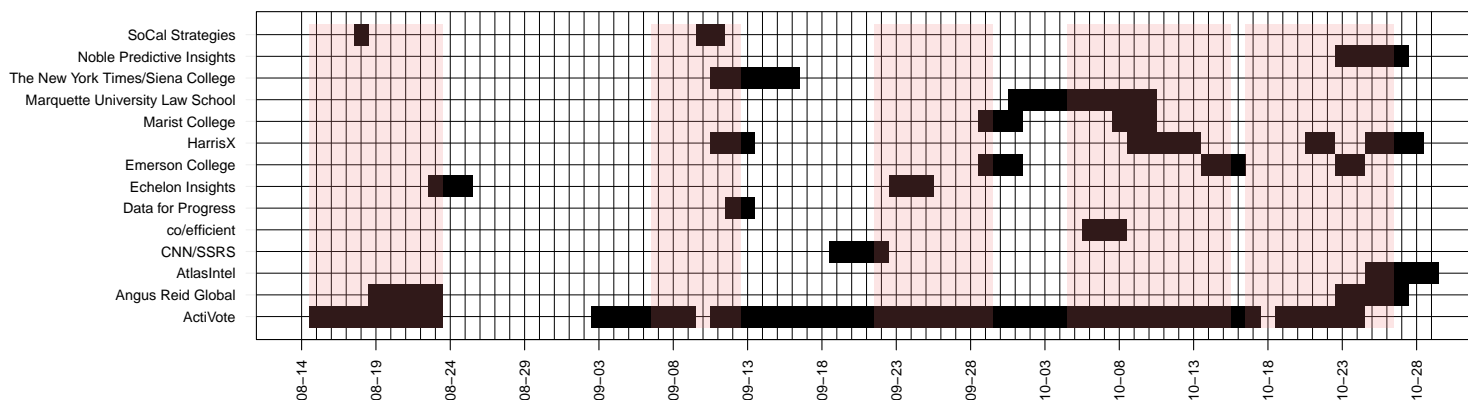


Figure C.18: Temporal coverage of reference pollsters for education-level crosstabs. The light red shading marks the digital fieldwork period for PoSSUM, during which these polls were combined into the average used to evaluate PoSSUM's results.

Temporal Coverage by Pollster for Crosstab: 2020 Vote

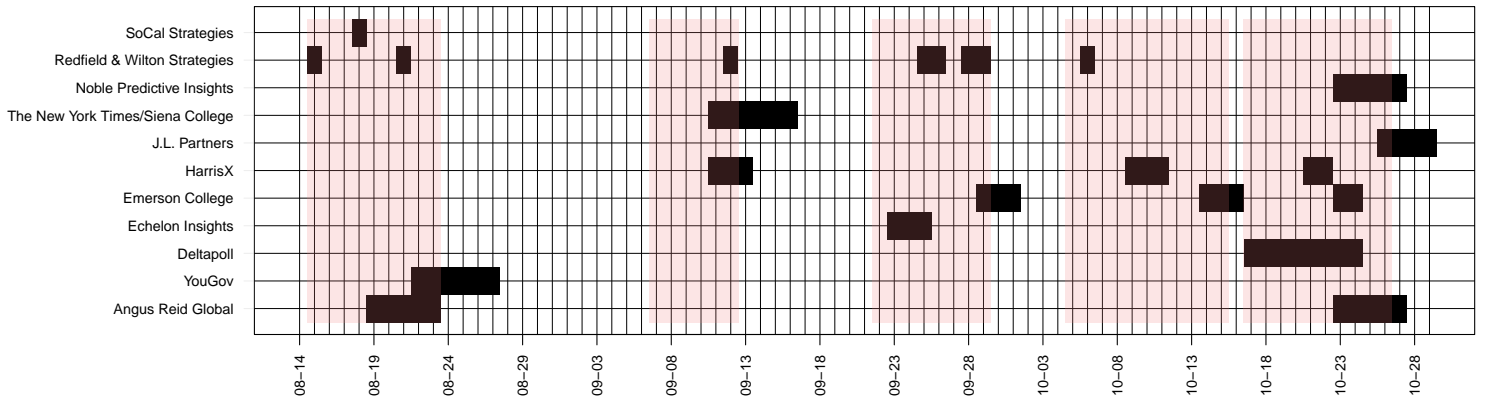


Figure C.19: Temporal coverage of reference pollsters for 2020 vote-level crosstabs. The light red shading marks the digital fieldwork period for PoSSUM, during which these polls were combined into the average used to evaluate PoSSUM's results.

C.2 Novel Learning, Human Alignment and Time-Sensitivity

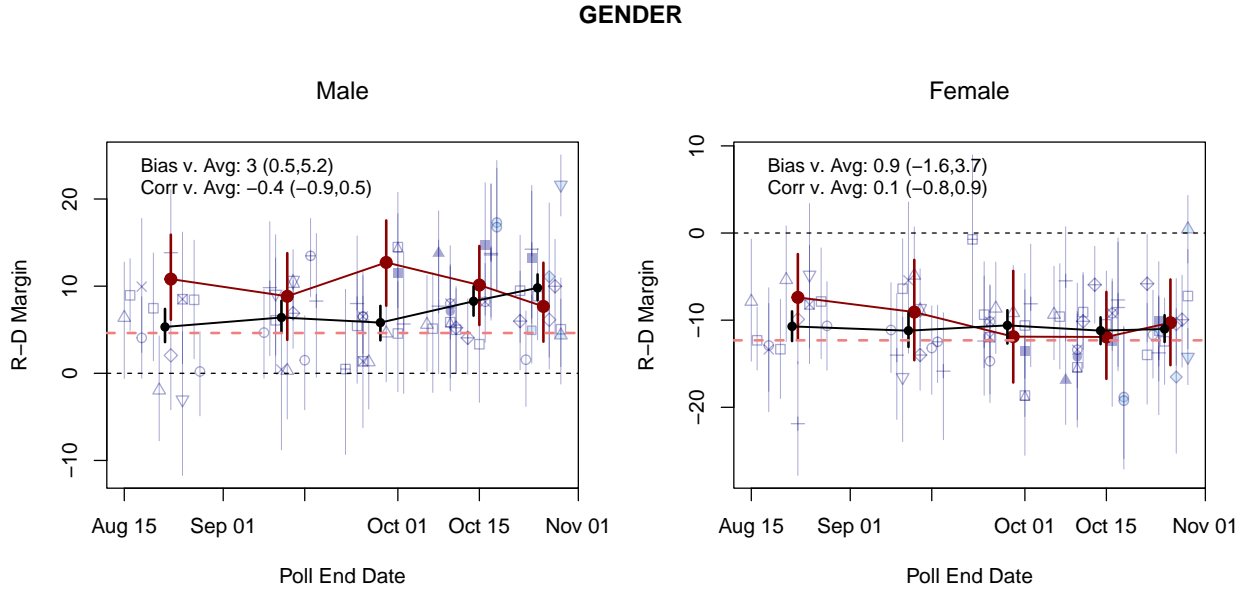


Figure C.20: Gender-level PoSSUM estimates over the course of the campaign, shown alongside individual polls overlapping PoSSUM's fieldwork periods, the aggregated polling average for each PoSSUM fieldwork window, and the reference preferences from 2020.

EDUCATION

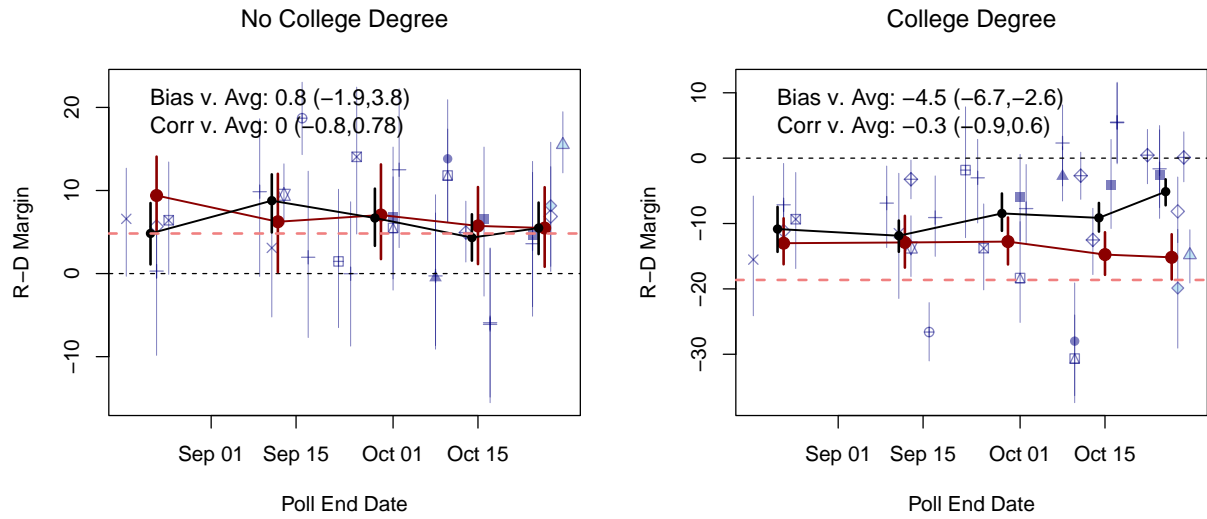


Figure C.21: Education-level PoSSUM estimates over the course of the campaign, shown alongside individual polls overlapping PoSSUM's fieldwork periods, the aggregated polling average for each PoSSUM fieldwork window, and the reference preferences from 2020.

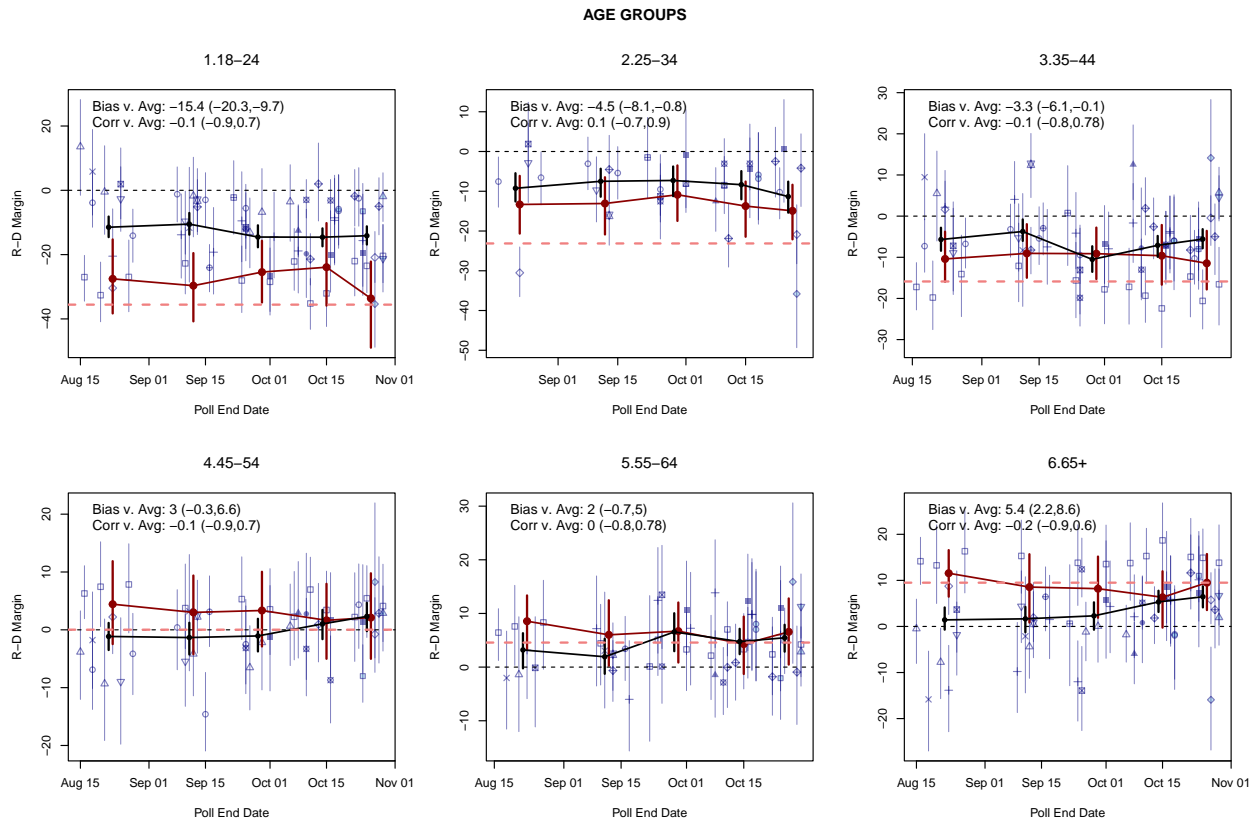


Figure C.22: Age-group-level PoSSUM estimates over the course of the campaign, shown alongside individual polls overlapping PoSSUM's fieldwork periods, the aggregated polling average for each PoSSUM fieldwork window, and the reference preferences from 2020.