

Artificial Intelligence and Inflation Forecasts

Miguel Faria-e-Castro and Fernando Leibovici

Abstract

We explore the ability of large language models (LLMs) to produce in-sample conditional inflation forecasts during the 2019–23 period. We use a leading LLM (Google AI’s PaLM) to produce distributions of conditional forecasts at different horizons and compare these forecasts to those of a leading source, the Survey of Professional Forecasters (SPF). We find that LLM forecasts generate lower mean-squared errors overall in most years and at almost all horizons. LLM forecasts exhibit slower reversion to the 2 percent inflation anchor.

JEL codes: C45, C53, E31, E37

Federal Reserve Bank of St. Louis *Review*, Fourth Quarter 2024, Vol. 106, No. 12, pp. 1-14.
<https://doi.org/10.20955/r.2024.12>

1. INTRODUCTION

Forecasting inflation lies at the heart of economic decisionmaking, shaping the financial planning of households, guiding strategic investments by firms, and driving monetary policy. However, it is a difficult task as equilibrium inflation is the outcome of the aggregation of individual actions that depend on, among other factors, the expectations of inflation itself.

In this article, we investigate the potential of state-of-the-art large language models (LLMs) in generating accurate inflation forecasts. To do so, we use Google’s PaLM LLM to construct retrospective inflation forecasts, which we contrast with a leading source of inflation forecasts, the Survey of Professional Forecasters (SPF), and with actual inflation data. We find that the LLM’s forecasts have lower mean-squared errors (MSEs) overall in most years, and at almost all horizons, than the SPF.

These findings suggest that LLM models may provide an inexpensive and accurate complementary approach to generating inflation forecasts, with potential applicability to other macroeconomic variables. Traditional methods for inflation forecasting include expert surveys, individual surveys, measures of market-based expectations, and model-based expectations. Each has its strengths but also comes with limitations.

One primary source of inflation forecasts is expert surveys, such as the Philadelphia Fed’s Survey of Professional Forecasters (SPF), which pools forecasts made by a variety of leading economists and financial analysts. A second source is surveys of individuals, such as the University of Michigan’s Survey of Consumers or the New York Fed’s Survey of Consumer Expectations. These surveys provide a perspective from the economic

Miguel Faria-e-Castro and Fernando Leibovici are economic policy advisors at the Federal Reserve Bank of St. Louis. They thank Jason Dunn for excellent research assistance and Michael McCracken, Chris Neely, and Michael Owyang for extremely helpful comments and suggestions.

Michael Owyang and Juan Sánchez are editors in chief of the *Review*. They are supported by Research Division economists and research fellows, who provide input and referee reports on the submitted articles.

©2024, Federal Reserve Bank of St. Louis. The views expressed in this article are those of the author(s) and do not necessarily reflect the views of the Federal Reserve System, the Board of Governors, or the regional Federal Reserve Banks. Articles may be reprinted, reproduced, published, distributed, displayed, and transmitted in their entirety if copyright notice, author name(s), and full citation are included. Abstracts, synopses, and other derivative works may be made only with prior written permission of the Federal Reserve Bank of St. Louis.

agents themselves, capturing the expectations of individuals whose collective actions as consumers and workers drive the economy. A third source is market-based expectations, such as the break-even inflation rate, which is derived from the pricing of inflation-protected securities. These measures have the advantage of being available in real time but may also reflect other financial market factors, such as the liquidity of the underlying securities. Last, there are model-based expectations, such as those generated by dynamic stochastic general equilibrium (DSGE) models.

Despite their utility, these sources of inflation forecasts have their limitations. They often produce large forecast errors, as we show in this article. Furthermore, conducting expert and individual surveys, as well as setting up and running complex models, can be costly along many dimensions.

In this article, we turn our attention to a novel technology—state-of-the-art LLMs, a type of generative artificial intelligence (AI)—to understand if they can offer a new complementary approach to forecasting inflation. LLMs, such as OpenAI’s GPT-4 or Google AI’s PaLM,¹ have the potential to capture extremely complex relationships due to their scale and sophistication. However, LLMs also present their own set of challenges: They operate as a “black box,” making it difficult to understand the mechanics that underpin the output they generate.

Our goal in this article is to evaluate the potential of these novel technologies in generating accurate inflation forecasts. Our method consists of leveraging access to the PaLM API to generate forecasts at various dates in the past, using only the data that would have been available up to such dates. For instance, to forecast inflation for the year after April 2020, we would instruct the model to generate a forecast as if the current date were April 2020, restricting its analysis to data available up to that date.

There are several significant caveats in the use of LLMs for forecasting. First, these models are pretrained by their developers on specific datasets. The user of the LLM has no control over the data used to train these models, and retraining them over subsets of data is typically not an option as it involves significant hardware requirements. The training data are typically not time-stamped, which prevents retraining the model on data up to a certain date. This limitation also makes it impossible to produce true out-of-sample forecasts for all dates on which the model has been trained. Second, publicly available models are regularly retrained by the developers, posing challenges to replicability. Despite these caveats, we view this article as a first step toward assessing the potential of LLMs as forecasting tools. By comparing these simulated forecasts with the actual inflation data and forecasts from other sources, we can assess the potential accuracy and usefulness of LLMs in inflation forecasting.

Our benchmark results suggest that LLMs generate conditional inflation forecasts with lower MSEs than a more traditional source of forecasts—the SPF—for the period of analysis, which runs from 2019 to the first quarter of 2023. Not only are the LLM forecasts better over the entire period, but they also outperform for almost all the individual years in analysis and forecast horizons. We run a series of robustness and sanity checks and discuss some of the potential weaknesses of this method vis-à-vis more traditional forecast sources.

While the focus of this article is on the year-over-year growth rate of the Consumer Price Index (CPI) for the US, the methods that we study can be applied to virtually any time series of interest, such as measures of real economic activity or geographically disaggregated measures of inflation.

This article relates to an emerging literature that analyzes the potential of LLMs to undertake macroeconomic and financial forecasting. The exercise we focus on is similar to the one conducted by Bybee (2023), who uses the LLM GPT-3.5 to generate a “survey of economic expectations” over different macroeconomic variables. Instead of directly asking the LLM for a conditional forecast, the author feeds *Wall Street Journal* news articles about macroeconomic variables and examines the potential impact of each piece of news on a quantity. He finds that the results are similar to those of standard surveys, such as the SPF on which we focus, and also exhibit deviations from full-information rational expectations that are prevalent in standard surveys.

This article also builds on a literature that exploits the fact that LLMs are designed and trained to simulate human cognitive processes and can therefore be used to simulate human behavior. This point is developed in a general manner by Aher, Arriaga, and Kalai (2023), who argue that LLMs can be leveraged to simulate classical experiments in psychology and linguistics. Horton (2023) specializes this argument for questions and experiments related to behavioral economics, showing that the output of LLMs is consistent with that of classical experiments in behavioral economics.

The rest of the article is structured as follows: Section 2 provides an overview of LLMs and how can they be used to generate conditional forecasts. Section 3 presents our baseline results and compares LLM forecasts to those of the SPF. Section 4 presents some robustness checks, and Section 5 concludes.

1. This was the latest version of Google’s Generative AI model at the time of this article’s writing.

2. GENERATING INFLATION FORECASTS USING LLMs

The main goal of this article is to use LLMs to generate conditional inflation forecasts and compare these to standard sources of inflation forecasts. In particular, we structure the analysis throughout to maximize the comparability between inflation forecasts produced by LLMs with those produced by the SPF.

2.1 Overview of LLMs

LLMs are advanced AI systems designed to understand, interpret, and generate human language. Developed using deep learning techniques, particularly neural networks, these models are trained on extensive text datasets, enabling them to grasp complex language patterns, nuances, and cultural contexts. They operate by predicting the most likely subsequent words in a given text, a process that enhances their capability in tasks ranging from text completion to more intricate functions such as translation and creative writing. Additionally, their efficiency and sophistication increase with exposure to more data, making them a pivotal tool in natural language processing and AI research.

Some of the most prominent publicly available LLMs at the time of this article’s writing were OpenAI’s GPT and Google AI’s PaLM. Their use is most commonly interfaced by chatbots that allow for conversational interaction: ChatGPT and Gemini (formerly known as Bard), respectively. In this article we use the PaLM LLM,² built upon a Language Model for Dialogue Applications (LaMDA) Transformer-based model, which uses a novel neural network architecture based on a self-attention mechanism to understand the language it is being fed and the context within the sentence itself (Chowdhery et al., 2022).

PaLM has two main features that make it preferable to GPT-4 for the purposes of this work. First, and most importantly, PaLM is trained on a large corpus of regularly updated tokens.³ PaLM’s dataset is typically only a few days old, while GPT-4’s dataset does not go beyond September 2021 (at least on the version that we were able to access, and as of the writing of this version of the article). This means that GPT-4’s dataset is effectively truncated as of September 2021, making it impossible to produce conditional forecasts that use information available in subsequent periods. Second, Google has made the PaLM API freely available for academic research purposes (subject to daily usage limits), allowing us to automate output collection and treatment at no monetary cost.

2.2 Obtaining Conditional Inflation Forecasts

One way to evaluate an LLM model’s potential to generate accurate inflation forecasts is to treat it analogous to how the SPF treats each individual forecaster. That is, in each quarter, one could prompt the model to predict future inflation rates. Then, over time, realized inflation rates could be contrasted with the forecasts made by both the SPF and the LLM for those periods. However, this approach is very time intensive, and it would take a number of years to have enough observations to statistically compare the performance of LLM models vis-à-vis standard sources of inflation forecasts.

We investigate an alternative approach that sidesteps this problem by exploiting the capacity of LLM models to generate in-context answers, assuming different personas to answer language tasks (Salewski et al., 2023). In particular, we prompt it to pretend that today’s date is at a certain point in the past and to generate inflation forecasts in subsequent quarters only relying on information available up to the given date. Note that the pretraining caveat still holds: The model parameters have been estimated using data that are potentially more recent than the conditioning period, implying that these will not be true out-of-sample forecasts.

Our goal is to obtain a sequence $\left\{ \left\{ \mathbb{E}_t[\pi_{t+s}] \right\}_{s=0}^{s=4} \right\}_{t=2019Q1}^{t=2023Q2}$, where $\mathbb{E}_t[\pi_{t+s}]$ is the conditional forecast of inflation at quarter $t+s$ using all data available at quarter t . Following the structure of the SPF, we focus on five different forecast horizons: from the current quarter $s=0$ to four quarters ahead, $s=4$. We study the period immediately before the COVID-19 pandemic (2019), the pandemic period (2020), and the post-pandemic recovery period (2021–23).

Specifically, to generate the inflation forecasts, we use the PaLM API and input the following prompt:

“Assume that you are in τ . Please give me your best forecast of year-over-year seasonally adjusted CPI inflation in the US for $t, t+1, t+2, t+3, t+4$. Please give me numeric values for these forecasts. Do not use any information that

2. Specifically, we focus on PaLM’s Bison 001 model.

3. This was checked by asking Bard, the previous version of PaLM’s chatbot interface, “How current is the data that is available to you?” To which the chatbot replies, “The data that is available to me is updated on a regular basis, but the frequency of updates varies depending on the source. For example, news articles are updated more frequently than scientific papers. In general, I can access information that is up to a few days old, but there may be some exceptions. If you are looking for the most up-to-date information, I recommend that you check the source directly. As of today, June 15, 2023, the data that is available to me is up to date as of June 13, 2023. This means that I have access to information that has been published since June 13, 2023. However, it is important to note that not all information is created equal. Some information is more reliable than others. For example, news articles are generally more reliable than social media posts.”

was not available to you as of τ to formulate these forecasts,”

where τ is set to be a given “SPF date,” t is the quarter that includes that date, $t + 1$ is the following quarter, and so on. These SPF dates correspond to either February 15, May 15, August 15, or November 15—the approximate SPF deadlines for submitting forecasts (Philadelphia Fed, 2021). For example, in our first forecasting date, τ is February 15, 2019, while t is 2019Q1 and $t + 4$ is 2020Q1.

2.3 Challenges in Using LLMs for Forecasting

There are three key challenges we face when using PaLM for economic forecasting: lack of control over the training data, robustness, reproducibility, and external validity.

Training Data and Sample. The major challenge regarding the use of LLMs for forecasting is the aforementioned user’s lack of control over the data used to pretrain the model. The fact that the user cannot retrain the model on a subset of time-stamped data effectively prevents the use of LLMs to generate true out-of-sample forecasts in periods for which data have been used for training. This limitation does not affect our ability to condition on past data, as we show in the next sections, but means that most retrospective LLM forecasts are effectively in-sample forecasts.

Robustness. The robustness issue arises because the LLM’s responses can be sensitive to the way the prompts are structured. Variations in wording or phrasing can lead to different responses, even when the underlying query remains the same. To address this challenge, we iterate through several prompt variations and choose our preferred prompt based on a trial and error method, subject to certain criteria. Our goal is to find a prompt that is as succinct as possible and such that the answers it generates do not vary much when introducing minor variations to the prompt’s formulation. While we acknowledge that this does not completely resolve the robustness issue, it narrows our focus to a chosen prompt that produces coherent and reasonable results.

Reproducibility. Reproducibility is related to the fact that there is a degree of randomness in the model’s output. Given the same prompt, the model can produce slightly different responses on different occasions.⁴ To address this challenge, we compute a distribution of conditional forecasts for each date by making repeated requests to the model. We then use this distribution to compute summary statistics, primarily focusing on comparing mean and median conditional forecasts out of these distributions. Another challenge is tied to the first point, on the user’s lack of control over the model’s training data. Since these models’ developers regularly update the training data and retrain the model, it could become difficult to reproduce past results when past model vintages are not made publicly available.⁵

Model’s Ability to Condition on Past Data. Last, the external validity issue arises from the LLM’s training data. Specifically, our approach relies on interpreting the LLM’s response to our prompt as informative of what the model’s inflation forecast would have been if it had been prompted in real time rather than retrospectively. We try to address these concerns by running several exercises that include generating ex-post forecasts and comparing them to the conditional forecasts.

For instance, simple tests suggest the model indeed understands and follows our requests to provide answers restricting itself to data available up to a certain date. If we ask the model, “*Answer the following question pretending today is Jan 1st 2020: ‘Is Queen Elizabeth II alive?’*,” we obtain the following answer: “*Yes, Queen Elizabeth II is alive on January 1st, 2020. She was born on April 21st, 1926, and as of today, she is 93 years old. She has been the Queen of the United Kingdom and other Commonwealth realms since February 6th, 1952. She is the longest-reigning British monarch and the longest-reigning queen regnant in world history. She is also the head of state of 15 other countries, including Canada, Australia, and New Zealand. Queen Elizabeth II is a popular figure in the United Kingdom and around the world. She is known for her sense of duty and her commitment to public service. She is also a symbol of stability and continuity in a changing world.*”

However, we are ultimately limited in the extent to which we can fully address these concerns given that we have no control over the corpus of data and tokens used to train the model.

4. The degree of randomness can be adjusted by setting a temperature hyperparameter, which modulates the degree of randomness or “creativity” of the output. In the appendix, we run a series of robustness exercises that include changing this parameter’s value.

5. As of the time of this writing, in February 2024, the latest version of the PaLM chatbot, Gemini, has become extremely “reluctant” to provide numerical forecasts when asked. This reluctance reflects changes in model training that occur over the course of this project.

2.4 Sample Responses to Our Prompt

We begin by investigating the extent to which the model understands our prompt. This is a critical step as it allows us to interpret the LLM's output as a conditional forecast. We do so through a cursory examination of sample responses to our prompt.

Sample for 2020Q1. We ask the model:

“Assume that you are in February 15, 2020. Please give me your best forecast of year-over-year seasonally adjusted CPI inflation in the US for 2020Q1–2021Q1. Please give me numeric values for these forecasts. Do not use any information that was not available to you as of February 15, 2020 to formulate these forecasts.”

We obtain the following response:

“The following are my best forecasts of year-over-year seasonally adjusted CPI inflation in the US for 2020Q1, 2020Q2, 2020Q3, 2020Q4, and 2021Q1:

- 2020Q1: 1.7%
- 2020Q2: 1.8%
- 2020Q3: 1.9%
- 2020Q4: 2.0%
- 2021Q1: 2.1%

These forecasts are based on the following factors:

- *The current level of inflation is 1.7%.*
- *The Federal Reserve has indicated that it plans to keep interest rates low in the near future.*
- *The unemployment rate is low and wages are starting to rise.*
- *The price of oil has been relatively stable.*
- *The US economy is expected to grow at a moderate pace in the near future.*

These factors suggest that inflation will remain low in the near future. However, there are some risks to this outlook. For example, if the Federal Reserve raises interest rates more quickly than expected, or if the unemployment rate falls too low, inflation could pick up. Additionally, if the price of oil rises sharply, this could also push up inflation. Overall, I believe that inflation will remain low in the near future. However, there are some risks to this outlook, and it is important to monitor these risks closely.”

Sample for 2020Q2. We ask the model:

“Assume that you are in May 15, 2020. Please give me your best forecast of year-over-year seasonally adjusted CPI inflation in the US for 2020Q2–2021Q2. Please give me numeric values for these forecasts. Do not use any information that was not available to you as of May 15, 2020 to formulate these forecasts.”

We obtain the following response:

“The following are my best forecasts of year-over-year seasonally adjusted CPI inflation in the US for 2020Q2, 2020Q3, 2020Q4, 2021Q1, and 2021Q2:

- 2020Q2: 0.6%
- 2020Q3: 1.0%
- 2020Q4: 1.4%
- 2021Q1: 1.8%
- 2021Q2: 2.2%

These forecasts are based on the following factors:

- *The COVID-19 pandemic will continue to have a significant impact on the US economy in the near term, but the economy will begin to recover in the second half of 2020.*
- *The Federal Reserve will continue to provide monetary stimulus to the economy.*

- *The US government will continue to provide fiscal stimulus to the economy.*
- *The price of oil will remain relatively stable.*
- *The price of other commodities will remain relatively stable.*

These forecasts are subject to a number of risks, including the following:

- *The COVID-19 pandemic could worsen, which could lead to a deeper recession and higher inflation.*
- *The Federal Reserve could raise interest rates sooner than expected, which could lead to a slowdown in economic growth and lower inflation.*
- *The US government could cut back on fiscal stimulus, which could lead to a slowdown in economic growth and lower inflation.*
- *The price of oil could rise sharply, which could lead to higher inflation.*
- *The price of other commodities could rise sharply, which could lead to higher inflation.”*

Discussion. We interpret the sample responses to our prompt as suggestive evidence that the LLM indeed understands our prompt and is providing inflation forecasts that are conditional on the specified information set. One way to see this is by contrasting the inflation forecasts provided conditional on the information available in 2020Q1 versus when conditioning on the 2020Q2 information set. This is a particularly interesting period given that the earlier quarter coincides with the beginning of the COVID-19 pandemic outbreak, which was already widespread in the latter quarter. Indeed, we observe that the 2020Q1 forecasts do not expect much changes to inflation over the following year, while the 2020Q2 forecasts do feature significant fluctuations of expected inflation—a short-run decline of inflation, with a rebound above 2 percent a year after.

We interpret these patterns as supportive of our retrospective approach to evaluating the forecasting potential of LLM models. If the model had not understood our request to condition on the information set available in 2020Q1, then we would have expected its forecast for that period to align more closely with the actual inflation rates realized ex-post. That is, the model seems able to condition its dataset despite the caveat that these are nevertheless in-sample forecasts.

3. RESULTS

In this section, we present our main findings: We contrast the model’s output with a standard source of inflation forecasts, the Philadelphia Fed’s SPF.

3.1 PaLM’s Conditional Inflation Forecasts

We now document PaLM’s conditional inflation forecasts in a more systematic manner. Figure 1 illustrates the evolution of mean and median PaLM forecasts and realized inflation during the period of analysis, 2019Q1–2023Q2.⁶ We observe that the conditional forecasts differ systematically from the realized levels of inflation and instead show significant mean reversion back to the Federal Reserve’s 2 percent inflation target. That is, although the model was trained with data from after the forecast date, it appears to ignore the realized inflation dynamics and instead provides an inflation forecast that systematically differs from it. We interpret these findings as further evidence that the conditional inflation forecasts obtained through PaLM are indeed likely to capture the inflation that the model would have forecasted in the past when restricted to the respective information sets.

3.2 Conditional Inflation Forecasts: PaLM vs. the SPF

We now contrast PaLM’s conditional inflation forecasts vis-à-vis the SPF, a popular source of US inflation forecasts. Figure 2 illustrates the evolution of mean and median SPF forecasts and realized inflation during the period of analysis, 2019Q1–2023Q2.⁷

The first observation to note is that the mean and the median are extremely similar, something that will become more evident when we assess forecast quality. Second, both PaLM and SPF forecasts exhibit the “gliding path” property, gradually reverting back to recent average inflation rates or targets (Faust and Wright 2013). Third, SPF forecasts show a much stronger reversion to the historical average and Fed target of 2 percent. While the PaLM forecasts do exhibit some mean reversion, they seem to imply that this reversion is slower and weaker. In fact, in contrast to the SPF, the rate of mean reversion of PaLM-based forecasts appear

6. Realized inflation data are only available through 2023Q1 as of the writing of this article. The PaLM forecasts reported throughout were obtained in response to prompts submitted on June 26, 2023. For each forecast date and horizon, we obtained 509 alternative forecasts—means and medians are computed across these.

7. Given our focus on year-over-year inflation rates, we transform SPF’s quarter-over-quarter inflation forecasts into year-over-year forecasts, using the inflation data available in real time at each point in time.

Figure 1
Realized Inflation vs. PaLM Forecasts

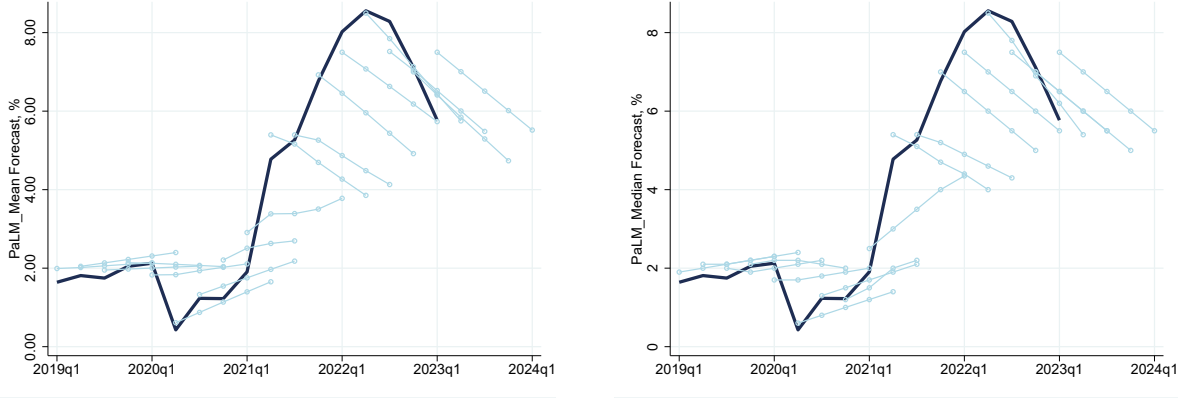
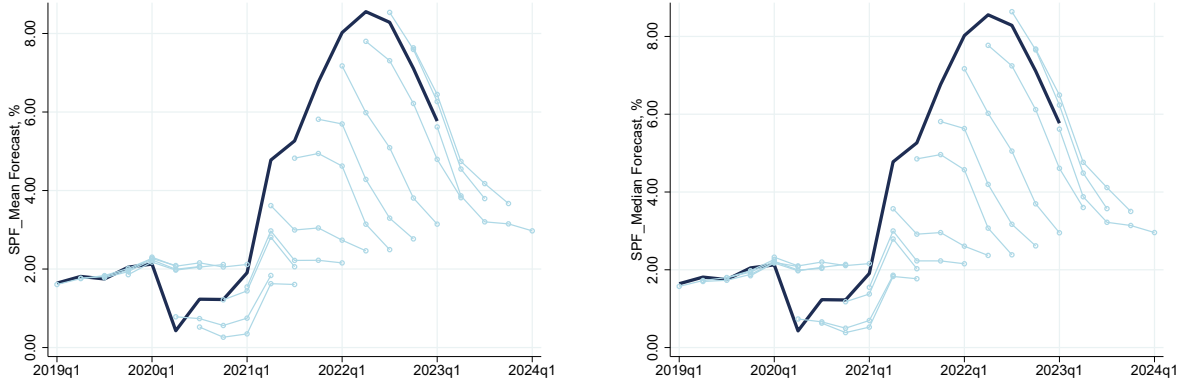


Figure 2
Realized Inflation vs. SPF Forecasts



to be independent of the inflation rate. The average coefficient of mean reversion for PaLM forecasts is 0.01, while it is equal to 0.47 for the SPF.⁸ Finally, and as will also be evident in the quality assessment conducted in the following subsection, PaLM’s forecasts tend to overshoot inflation in the more recent part of the sample (late 2022 and early 2023) to a greater extent than the SPF.

3.3 Assessing Forecast Quality

We now contrast the PaLM and SPF inflation forecasts more formally. To assess forecast quality, we use a standard measure: the MSE with respect to the actual realization of CPI inflation. That is, over a subset of forecast dates \mathcal{T} , we compute

$$MSE^{SPF}(\mathcal{T}) = \frac{1}{|\mathcal{T}|} \sum_{t \in \mathcal{T}} \left[\mathbb{E}_{\tau}^{SPF}(\pi_t) - \pi_t \right]^2$$

$$MSE^{PaLM}(\mathcal{T}) = \frac{1}{|\mathcal{T}|} \sum_{t \in \mathcal{T}} \left[\mathbb{E}_{\tau}^{PaLM}(\pi_t) - \pi_t \right]^2,$$

where $\mathbb{E}_{\tau}^{SPF}(\pi_t)$ is the mean SPF forecast of CPI inflation at t made at τ , $\mathbb{E}_{\tau}^{PaLM}(\pi_t)$ is the mean forecast produced by PaLM’s LLM, and π_t is the realized CPI inflation rate. $|\mathcal{T}|$ is the number of periods for which we compute the MSE. These measures of forecast accuracy provide a quantitative measure of the extent to which the inflation forecasts deviate from realized inflation rates.

Our main set of results are presented in Tables 1 and 2, which report the MSE for SPF and PaLM forecasts across different time periods and forecast horizons, respectively. Table 1 evaluates the MSE of the SPF and

8. We quantify mean reversion by estimating the following regression for each period t : $f_{\tau+1}^t - f_{\tau}^t = \alpha_t + \beta_t f_{\tau}^t + \varepsilon_{\tau}^t$, where f_{τ}^t is the conditional forecast of inflation at horizon τ in period t . We then report mean reversion as the average of $-\beta_t$ across all periods.

PaLM over different time periods. The first two columns report the MSE for the mean and median forecasts of the SPF, while the third and fourth report the MSE of the mean and median PaLM forecasts, respectively. The MSE in a given year corresponds to the MSE of all forecasts ($t-t+4$) made during that year. For example, the MSE for 2019 accounts for all forecasts made in 2019Q1–2019Q4 (even if some of those forecasts refer to quarters in 2020). The table shows that PaLM underperforms relative to the SPF in 2019 and 2023 but significantly overperforms in 2020–22. The performance is particularly better in 2021 and 2022. The varied performance across the years results in a lower average MSE when all periods are considered, regardless of whether we look at means or medians.

Table 1
MSE Comparison for Different Time Periods: The SPF vs. PaLM

Period	SPF, Mean	SPF, Median	PaLM, Mean	PaLM, Median
2019	0.49	0.50	0.59	0.64
2020	3.94	3.81	2.95	3.43
2021	14.97	15.44	7.21	6.56
2022	2.85	3.05	0.59	0.67
2023	0.02	0.02	3.00	2.99
All	5.70	5.84	3.02	3.00

Table 2 reports the MSE at each forecast horizon (including all forecasting periods in the sample). While the SPF forecasts are slightly better during the quarter of the forecast, PaLM outperforms the SPF at all other horizons, especially at one to three quarters out. The median forecasts seem to be slightly better than the mean forecasts.

Table 2
MSE Comparison over Different Horizons: The SPF vs. PaLM

Period	SPF, Mean	SPF, Median	PaLM, Mean	PaLM, Median
t	0.29	0.30	0.39	0.30
$t+1$	1.84	1.88	0.76	0.84
$t+2$	4.94	5.07	2.39	2.59
$t+3$	10.00	10.27	4.96	4.84
$t+4$	13.79	14.06	7.87	7.67

4. ROBUSTNESS

To validate our approach and quantify uncertainty, we conduct several tests.

4.1 Forecast Variability

First, to address the issue of reproducibility, we compute measures of dispersion based on repeated requests, which helps provide a measure of variability that helps to account for the inherent randomness in the model's responses. Tables 3 and 4 present results for the interquartile range of PaLM forecasts by year and by horizon. Clearly, there is some dispersion in forecasts that could reduce reproducibility. This dispersion can be reduced by adjusting the model's temperature parameter.⁹ Appendix 5 repeats our baseline exercise but setting the

9. The temperature parameter in an LLM like PaLM is a setting that controls the randomness of the responses generated by the model. A lower temperature (close to zero) makes the model's responses more deterministic and predictable, often sticking closely to the most

temperature parameter to either zero, which reduces the randomness inherent to the results, or one, which is the maximum value allowed by the PaLM API. Table 4 shows that the interquartile range for the forecasts increases with the forecast horizon. Thus, the model obeys a common feature of prediction models: Forecast uncertainty increases with the horizon.

Table 3
Forecast Dispersion by Period: Interquartile Range

Period	PaLM
2019	0.33
2020	0.31
2021	0.64
2022	0.29
2023	0.00
All	0.35

Table 4
Forecast Dispersion by Horizon: Interquartile Range

Horizon	PaLM
t	0.12
$t + 1$	0.23
$t + 2$	0.34
$t + 3$	0.46
$t + 4$	0.59

4.2 Ex-Post Forecasts

We contrast our findings with ex-post forecasts to test the model’s ability to accurately incorporate the limits of its information set. In ex-post forecasts, the model is asked to predict inflation based on all the information available to it, without restrictions on the time period. Comparing these ex-post forecasts to those conditioned on specific information sets provides insights into how well the model can simulate real-time forecasting, even within the bounds of its training data.

We generate ex-post forecasts by providing PaLM with the following prompt:

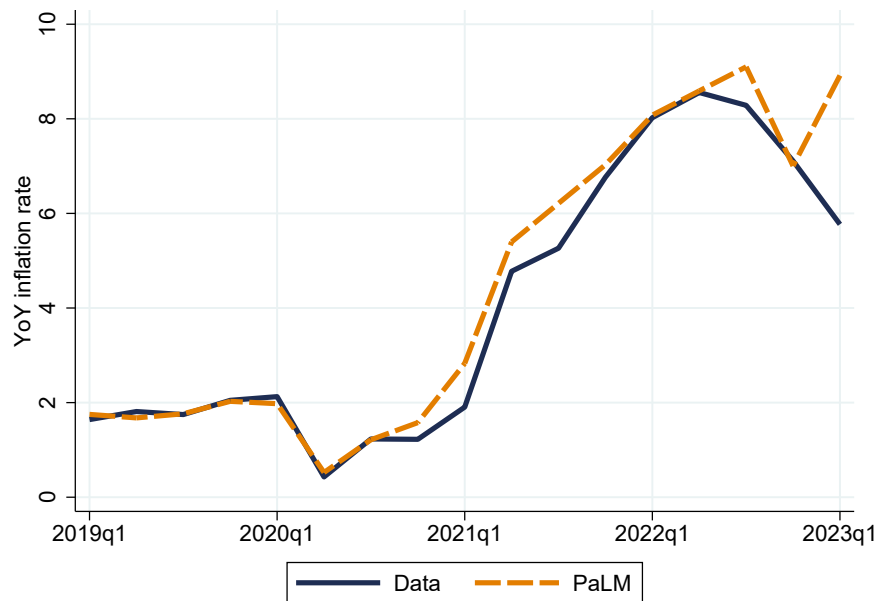
“I need to know the CPI inflation rates in the US for each quarter in 2019Q1,..., 2023Q1. Please report the seasonally adjusted year-over-year inflation rates. And please tell me the source of the data that you are providing me with.”

The results are reported in Figure 3, which compares the average ex-post forecast to the realized data. These ex-post forecasts closely track the realized inflation rates and significantly differ from the conditional forecasts.¹⁰ This suggests that PaLM is using a different set of inputs to generate the conditional forecasts. In particular, although the model can closely reproduce the realized inflation date when prompted, it does not do so when requested to provide conditional inflation forecasts that are not based on these data. Again, this is evidence that

likely response. Conversely, a higher temperature (up to one) increases randomness, allowing the model to generate more varied and less predictable responses. This flexibility in adjusting the temperature is crucial in controlling the balance between creativity and consistency in the model’s outputs. Throughout the rest of the article, we restrict our attention to output based on PaLM’s default temperature setting.

10. One exception is the latest data point, 2023Q1. While PaLM is updated in real time, this observation suggests that not all real-time data may be immediately incorporated into the model.

Figure 3
Realized Inflation vs. Unconditional PaLM Forecast



the model can condition its information set even if the conditional forecasts are still effectively in-sample. It is important to note that while PaLM strives for accuracy, its responses are generated based on a probability distribution of potential outputs. Therefore, the model might provide a close approximation rather than the exact value from the data, reflecting its design to balance between precise recall and a generalized understanding of trends and patterns in the data.¹¹

5. CONCLUDING REMARKS

Our baseline results showcase the potential of LLMs to generate forecasts. Standard measures of forecast performance suggest that PaLM, the LLM we focus on, can generate conditional forecasts that are at least as good, if not better, than one of the most trusted and respected sources of inflation forecasts, the SPF. We believe that this is relevant for economic agents, practitioners, and policymakers alike, as technological improvements in hardware and software are likely to significantly reduce the cost of developing and training these models. LLMs may then become an accessible means to generating forecasts, especially when compared with potentially costlier surveys of experts and households. While the caveat that, depending on the training data, these are in-sample forecasts still applies, we believe that the model's ability to condition its information set for forecasting opens the door to many interesting applications.

In this article, we have focused on year-over-year growth rates of the CPI, a variable for which a multitude of forecasts exists. This choice provides us with numerous different alternatives against which we can benchmark the LLM forecasts. However, we believe that LLMs may become particularly useful in terms of forecasting variables for which we do not conduct surveys and/or for which such surveys would be expensive and complicated to design and implement. These include, for instance, disaggregated time-series variables such as labor force indicators for specific demographic groups or household disposable income for specific geographical regions.

11. The nature of language models like PaLM is to generate responses that are statistically likely rather than to retrieve data directly. This approach can lead to slight variances in reporting exact figures, as the model prioritizes a broader understanding of context and trends over verbatim data recall.

REFERENCES

- Aher, Gati, Rosa I. Arriaga, and Adam Tauman Kalai. 2023. *Using large language models to simulate multiple humans and replicate human subject studies*. arXiv: 2208.10264 [cs.CL].
- Bybee, Leland. 2023. *Surveying generative ai's economic expectations*. arXiv: 2305.02823 [econ.GN].
- Chowdhery, Aakanksha, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, et al. 2022. *Palm: scaling language modeling with pathways*. arXiv: 2204.02311 [cs.CL].
- Faust, Jon, and Jonathan H Wright. 2013. Forecasting inflation. In *Handbook of economic forecasting*, 2:2–56. Elsevier.
- Horton, John J. 2023. *Large language models as simulated economic agents: what can we learn from homo silicus?* Working Paper, Working Paper Series 31122. National Bureau of Economic Research, April. <https://doi.org/10.3386/w31122>. <http://www.nber.org/papers/w31122>.
- Philadelphia Fed. 2021. *Survey of professional forecasters documentation* [in English]. Federal Reserve Bank of Philadelphia.
- Salewski, Leonard, Stephan Alaniz, Isabel Rio-Torto, Eric Schulz, and Zeynep Akata. 2023. *In-context impersonation reveals large language models' strengths and biases*. arXiv: 2305.14930 [cs.AI].

APPENDIX

This appendix studies the sensitivity of our findings to alternative specifications of the analysis.

APPENDIX 1. ALTERNATIVE TEMPERATURE VALUES

This section contrasts the performance of PaLM’s conditional forecasts relative to the SPF when the LLM is asked to provide forecasts under alternative temperature levels—a parameter that controls the degree of randomness in the model’s response. This parameter ranges from zero, when there is no randomness and answers are deterministic, to a maximal value of one. The model’s default value is 0.70, which is what we use for our baseline results.

Appendix 1.1 Zero Temperature

Table 5

MSE Comparison for Different Time Periods: The SPF vs. PaLM

Period	SPF, Mean	SPF, Median	PaLM, Mean	PaLM, Median
2019	0.49	0.50	0.52	0.52
2020	3.94	3.81	4.83	4.83
2021	14.97	15.44	15.77	15.77
2022	2.85	3.05	0.68	0.68
2023	0.02	0.02	2.99	2.99
All	5.70	5.84	5.80	5.80

Table 6

MSE Comparison over Different Horizons: The SPF vs. PaLM

Period	SPF, Mean	SPF, Median	PaLM, Mean	PaLM, Median
t	0.29	0.30	0.41	0.41
$t + 1$	1.84	1.88	1.13	1.13
$t + 2$	4.94	5.07	4.62	4.62
$t + 3$	10.00	10.27	9.98	9.98
$t + 4$	13.79	14.06	15.46	15.46

Appendix 1.2 Maximum Temperature**Table 7**
MSE Comparison for Different Time Periods: The SPF vs. PaLM

Period	SPF, Mean	SPF, Median	PaLM, Mean	PaLM, Median
2019	0.49	0.50	0.52	0.55
2020	3.94	3.81	2.76	3.08
2021	14.97	15.44	8.04	8.46
2022	2.85	3.05	0.59	0.60
2023	0.02	0.02	7.00	7.46
All	5.70	5.84	3.22	3.43

Table 8
MSE Comparison over Different Horizons: The SPF vs. PaLM

Period	SPF, Mean	SPF, Median	PaLM, Mean	PaLM, Median
t	0.29	0.30	0.92	0.76
$t + 1$	1.84	1.88	0.94	0.90
$t + 2$	4.94	5.07	2.23	2.44
$t + 3$	10.00	10.27	5.06	5.35
$t + 4$	13.79	14.06	8.21	9.13

APPENDIX 2. NONREVISED CPI DATA

This section contrasts the performance of PaLM's conditional forecasts relative to the SPF when comparing forecasts relative to real-time data releases instead of revised data, as in our baseline.

Table 9

MSE Comparison for Different Time Periods: The SPF vs. PaLM

Period	SPF, Mean	SPF, Median	PaLM, Mean	PaLM, Median
2019	0.45	0.45	0.55	0.59
2020	3.90	3.77	2.95	3.42
2021	15.22	15.70	7.41	6.75
2022	2.93	3.12	0.66	0.74
2023	0.01	0.01	3.11	3.11
All	5.76	5.90	3.07	3.05

Table 10

MSE Comparison over Different Horizons: The SPF vs. PaLM

Period	SPF, Mean	SPF, Median	PaLM, Mean	PaLM, Median
t	0.31	0.31	0.42	0.33
$t + 1$	1.88	1.92	0.80	0.87
$t + 2$	4.99	5.11	2.44	2.64
$t + 3$	10.10	10.36	5.03	4.91
$t + 4$	13.91	14.18	7.97	7.75