

Supporting Online Material
The Parable of Google Flu: Traps in Big Data Analysis

Authors: David Lazer^{1,2*}, Ryan Kennedy^{1,2,3}, Gary King², Alessandro Vespignani¹

Affiliations:

¹ Northeastern University, Boston, MA

² Harvard University, Cambridge, MA.

³ University of Houston, Houston, TX.

*Corresponding author. E-mail: d.lazer@neu.edu.

Contents:

1. Comparison of Time Series and GFT Results Across HHS Regions.
2. Autocorrelation and Partial Autocorrelation analysis of CDC % ILI data and GFT errors.
3. Tracking Google's Search Algorithm: An Analysis of Google's Official Search Blog
4. Search Dynamics – Which Search Terms Likely Led to GFT's Inaccuracy and the Problem of Replicating GFT's Results.
5. A Note on GFT's 2008 Model
6. Lag Models Using CDC Data 3 and 4 Weeks Out
7. Full Results and Sensitivity Tests
8. Comparison With Nonlinear Models

1. Comparison of Times Series and GFT Results Across HHS Regions

As we noted in the main document, national-level flu trends can be well-approximated using lagged dependent variables and variables for seasonality. These approximations can be further improved by combining GFT's estimates with a lagged dependent variable, 2 and 3 week lags of the difference between GFT and the CDC data, and the seasonality variables.

This section demonstrates that these results hold up in a region-by-region analysis. Table S1 compares errors across the entire period from 2003 to 2013 for the GFT model, the time series model (both with 07/03/2011, when GFT went off the rails, and with 09/06/2009, when the new GFT algorithm went live, as the out-of-sample start date). Table S2 compares the results in only the out-of-sample data for the post-09/06/2009 model. Table S3 compares the results in only the out-of-sample data for the post-07/03/2011 model. In both of these latter two tables, the comparison is against GFT after these respective dates.

The use of regional comparisons is closer to the method used by the original GFT, which was developed primarily in relation to what was then nine regions. The results are generally consistent with the national results. GFT's errors, however, are not uniformly distributed across regions. For example, in the post-2011 period, GFT does remarkably well in regions 7 and 8, while performing very poorly in regions 6 and 10. Nevertheless, in almost every case, GFT is out-performed by the basic time series predictions and the combined model. While not discussed in the article, it is interesting to note that GFT started to estimate high regularly after the summer 2011 change in geocoding, which was aimed at providing more accurate geographic results for searches (see section 3 of SM). It is possible that this contributed to the algorithm's later increase in error (i.e. greater accuracy in geographic search attribution threw off an algorithm designed on the less accurate attribution).

Replication code and data for these results are available in the folder labeled SOM1 in the replication files. The replication data is labeled `ParableOfGFT(SOM1Replication).dta` and the code is `SOMpt1(Replication Code).do`. Instructions for replication are included in the code file. The Stata statistical program was used for calculations.

Table S1: Regional Comparison of Models

	Google Flu Trends	CDC Lag and Seasonality Model (Post-07/03/2011 Out-of-Sample)	Google Flu Trends Plus Lag and Seasonality Corrections (Post-07/03/2011 Out-of-Sample)	CDC Lag and Seasonality Model (Post-09/06/2009 Out-of-Sample)	Google Flu Trends Plus Lag and Seasonality Corrections (Post-09/06/2009 Out-of-Sample)
Region 1 (CT, ME, MA, NH, RI, VT)	RMSE = 0.772 MAE = 0.276	RMSE = 0.531 MAE = 0.297	RMSE = 0.532 MAE = 0.216	RMSE = 0.581 MAE = 0.292	RMSE = 0.538 MAE = 0.219
Region 2 (NJ, NY)	RMSE = 0.876 MAE = 0.527	RMSE = 0.615 MAE = 0.392	RMSE = 0.474 MAE = 0.309	RMSE = 0.640 MAE = 0.397	RMSE = 0.487 MAE = 0.313
Region 3 (DE, DC, MD, PA, VA, WV)	RMSE = 0.832 MAE = 0.556	RMSE = 0.682 MAE = 0.439	RMSE = 0.513 MAE = 0.337	RMSE = 0.717 MAE = 0.448	RMSE = 0.526 MAE = 0.351
Region 4 (AL, FL, GA, KY, MS, NC, SC, TN)	RMSE = 0.695 MAE = 0.380	RMSE = 0.491 MAE = 0.310	RMSE = 0.353 MAE = 0.203	RMSE = 0.500 MAE = 0.316	RMSE = 0.358 MAE = 0.207
Region 5 (IL, IN, MI, MN, OH, WI)	RMSE = 0.738 MAE = 0.379	RMSE = 0.576 MAE = 0.342	RMSE = 0.390 MAE = 0.214	RMSE = 0.599 MAE = 0.344	RMSE = 0.400 MAE = 0.221
Region 6 (AR, LA, NM, OK, TX)	RMSE = 1.280 MAE = 0.732	RMSE = 0.881 MAE = 0.570	RMSE = 0.728 MAE = 0.441	RMSE = 0.913 MAE = 0.582	RMSE = 0.760 MAE = 0.445
Region 7 (IA, KS, MO, NE)	RMSE = 0.894 MAE = 0.508	RMSE = 0.774 MAE = 0.442	RMSE = 0.492 MAE = 0.281	RMSE = 0.812 MAE = 0.453	RMSE = 0.507 MAE = 0.289
Region 8 (CO, MT, ND, SD, UT, WY)	RMSE = 0.610 MAE = 0.342	RMSE = 0.555 MAE = 0.319	RMSE = 0.434 MAE = 0.248	RMSE = 0.595 MAE = 0.323	RMSE = 0.452 MAE = 0.250
Region 9 (AZ, CA, HI, NV)	RMSE = 0.971 MAE = 0.623	RMSE = 0.659 MAE = 0.435	RMSE = 0.525 MAE = 0.344	RMSE = 0.672 MAE = 0.445	RMSE = 0.532 MAE = 0.364
Region 10 (AK, ID, OR, WA)	RMSE = 1.050 MAE = 0.610	RMSE = 0.818 MAE = 0.520	RMSE = 0.710 MAE = 0.450	RMSE = 0.839 MAE = 0.538	RMSE = 0.717 MAE = 0.473

RMSE = Root Mean Squared Error

MAE = Mean Absolute Error

Note: The lagged CDC model is specified as

$$flu_t = \alpha + \beta_1 flu_{t-2} + \beta_2 flu_{t-3} + \beta_3 flu_{t-4} + \sum_{i=1}^{52} \gamma_i week_{it}$$

where flu is the CDC estimate of percent doctors' visits for ILI, $week$ is a binary variable indicating week of observation, and β and γ are estimated regression coefficients. The combination of GFT and CDC data is specified as

$$flu_t = \alpha + \beta_1 gflu_t + \beta_2 flu_{t-2} + \beta_3 (gflu_{t-2} - flu_{t-2}) + \beta_4 (gflu_{t-3} - flu_{t-3}) + \sum_{i=1}^{52} \gamma_i week_{it}$$

where $gflu$ is the GFT estimate. (See Section 7 of SM for details.)

Table S2: Regional Comparison of Models Post-09/06/2009

	Google Flu Trends	CDC Lag and Seasonality Model (Post-09/06/2009 Out-of-Sample)	Google Flu Trends Plus Lag and Seasonality Corrections (Post-09/06/2009 Out-of-Sample)
Region 1 (CT, ME, MA, NH, RI, VT)	RMSE = 1.152 MAE = 0.375	RMSE = 0.741 MAE = 0.325	RMSE = 0.788 MAE = 0.291
Region 2 (NJ, NY)	RMSE = 1.113 MAE = 0.617	RMSE = 0.727 MAE = 0.417	RMSE = 0.569 MAE = 0.346
Region 3 (DE, DC, MD, PA, VA, WV)	RMSE = 1.029 MAE = 0.689	RMSE = 0.858 MAE = 0.492	RMSE = 0.603 MAE = 0.381
Region 4 (AL, FL, GA, KY, MS, NC, SC, TN)	RMSE = 0.960 MAE = 0.590	RMSE = 0.477 MAE = 0.323	RMSE = 0.392 MAE = 0.231
Region 5 (IL, IN, MI, MN, OH, WI)	RMSE = 0.988 MAE = 0.499	RMSE = 0.666 MAE = 0.385	RMSE = 0.520 MAE = 0.275
Region 6 (AR, LA, NM, OK, TX)	RMSE = 1.584 MAE = 0.824	RMSE = 0.993 MAE = 0.591	RMSE = 0.904 MAE = 0.479
Region 7 (IA, KS, MO, NE)	RMSE = 1.033 MAE = 0.611	RMSE = 0.941 MAE = 0.522	RMSE = 0.651 MAE = 0.355
Region 8 (CO, MT, ND, SD, UT, WY)	RMSE = 0.740 MAE = 0.373	RMSE = 0.728 MAE = 0.357	RMSE = 0.546 MAE = 0.266
Region 9 (AZ, CA, HI, NV)	RMSE = 1.308 MAE = 0.872	RMSE = 0.630 MAE = 0.408	RMSE = 0.592 MAE = 0.385
Region 10 (AK, ID, OR, WA)	RMSE = 1.256 MAE = 0.729	RMSE = 0.748 MAE = 0.473	RMSE = 0.668 MAE = 0.437

RMSE = Root Mean Squared Error

MAE = Mean Absolute Error

Note: The lagged CDC model is specified as

$$flu_t = \alpha + \beta_1 flu_{t-2} + \beta_2 flu_{t-3} + \beta_3 flu_{t-4} + \sum_{i=1}^{52} \gamma_i week_{it}$$

where flu is the CDC estimate of percent doctors' visits for ILI, $week$ is a binary variable indicating week of observation, and β and γ are estimated regression coefficients. The combination of GFT and CDC data is specified as

$flu_t = \alpha + \beta_1 gflu_t + \beta_2 flu_{t-2} + \beta_3 (gflu_{t-2} - flu_{t-2}) + \beta_4 (gflu_{t-3} - flu_{t-3}) + \sum_{i=1}^{52} \gamma_i week_{it}$
 where $gflu$ is the GFT estimate. (See Section 7 of SM for details.)

Table S3: Regional Comparison of Models Post-07/03/2011

	Google Flu Trends	CDC Lag and Seasonality Model (Post-07/03/2011 Out-of-Sample)	Google Flu Trends Plus Lag and Seasonality Corrections (Post-07/03/2011 Out-of-Sample)
Region 1 (CT, ME, MA, NH, RI, VT)	RMSE = 1.526 MAE = 0.499	RMSE = 0.345 MAE = 0.216	RMSE = 1.004 MAE = 0.351
Region 2 (NJ, NY)	RMSE = 1.322 MAE = 0.623	RMSE = 0.483 MAE = 0.312	RMSE = 0.552 MAE = 0.324
Region 3 (DE, DC, MD, PA, VA, WV)	RMSE = 1.320 MAE = 0.951	RMSE = 0.594 MAE = 0.390	RMSE = 0.611 MAE = 0.372
Region 4 (AL, FL, GA, KY, MS, NC, SC, TN)	RMSE = 1.117 MAE = 0.610	RMSE = 0.452 MAE = 0.308	RMSE = 0.435 MAE = 0.253
Region 5 (IL, IN, MI, MN, OH, WI)	RMSE = 1.284 MAE = 0.658	RMSE = 0.459 MAE = 0.305	RMSE = 0.606 MAE = 0.301
Region 6 (AR, LA, NM, OK, TX)	RMSE = 2.013 MAE = 1.091	RMSE = 0.737 MAE = 0.486	RMSE = 0.894 MAE = 0.484
Region 7 (IA, KS, MO, NE)	RMSE = 0.524 MAE = 0.363	RMSE = 0.628 MAE = 0.398	RMSE = 0.714 MAE = 0.344
Region 8 (CO, MT, ND, SD, UT, WY)	RMSE = 0.550 MAE = 0.329	RMSE = 0.358 MAE = 0.242	RMSE = 0.512 MAE = 0.233
Region 9 (AZ, CA, HI, NV)	RMSE = 1.441 MAE = 0.728	RMSE = 0.539 MAE = 0.352	RMSE = 0.659 MAE = 0.375
Region 10 (AK, ID, OR, WA)	RMSE = 1.613 MAE = 0.977	RMSE = 0.464 MAE = 0.354	RMSE = 0.744 MAE = 0.437

RMSE = Root Mean Squared Error

MAE = Mean Absolute Error

Note: The lagged CDC model is specified as

$$flu_t = \alpha + \beta_1 flu_{t-2} + \beta_2 flu_{t-3} + \beta_3 flu_{t-4} + \sum_{i=1}^{52} \gamma_i week_{it}$$

where flu is the CDC estimate of percent doctors' visits for ILI, $week$ is a binary variable indicating week of observation, and β and γ are estimated regression coefficients. The combination of GFT and CDC data is specified as

$$flu_t = \alpha + \beta_1 gflu_t + \beta_2 flu_{t-2} + \beta_3 (gflu_{t-2} - flu_{t-2}) + \beta_4 (gflu_{t-3} - flu_{t-3}) + \sum_{i=1}^{52} \gamma_i week_{it}$$

where $gflu$ is the GFT estimate. (See Section 7 of SM for details.)

2. Autocorrelation and Partial Autocorrelation analysis of CDC % ILI data and GFT errors

As noted in the main text, the signs of autocorrelation in the CDC's data and in GFT's errors are apparent from simple examination of the time series charts. A more formal method for picturing this are correlograms and partial correlograms. Fig S1 shows the autocorrelations and 95% confidence intervals (grey shaded region) for the raw CDC percent visits for influenza-like illness (ILI) data. These results are very highly autocorrelated up to nearly 8 lags. Fig S2 shows the partial correlogram, which shows that the partial autocorrelation extends for about two lags.

Perhaps the most interesting plot here is Fig S3, however, which shows that the GFT model produces highly autocorrelated errors. This means that we can predict quite accurately how GFT will miss by in time t by looking at how much it misses by in times $t-1$, $t-2$, etc. This is the intuition that leads us to using the lagged error along with GFT's data to produce our most accurate model.

Replication of data for this figure is in the Manuscript file in the replication file, labelled `ParableOfGFT(Replication).dta`. Code is located in the SOM2 folder with the code `SOM2(Replication Code).do`. Calculations done in Stata.

Fig S1: Correlogram of CDC % ILI Data

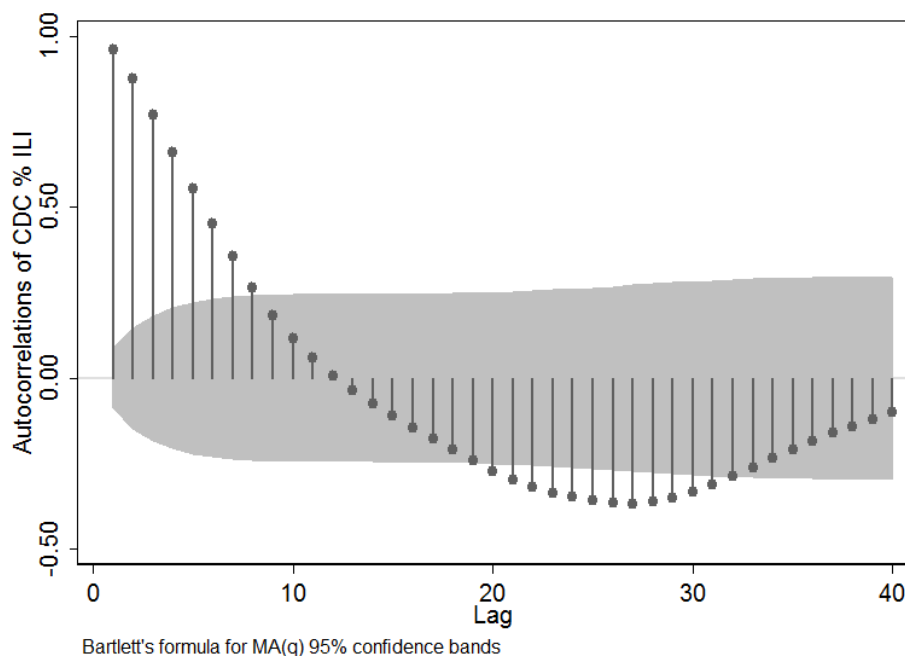


Fig S2: Partial Correlogram of CDC % ILI Data

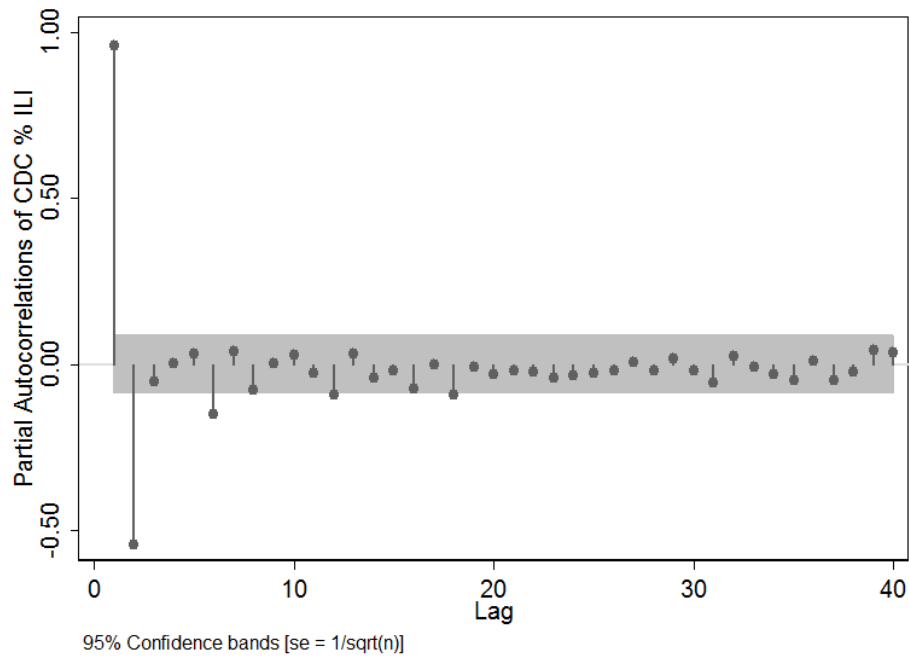
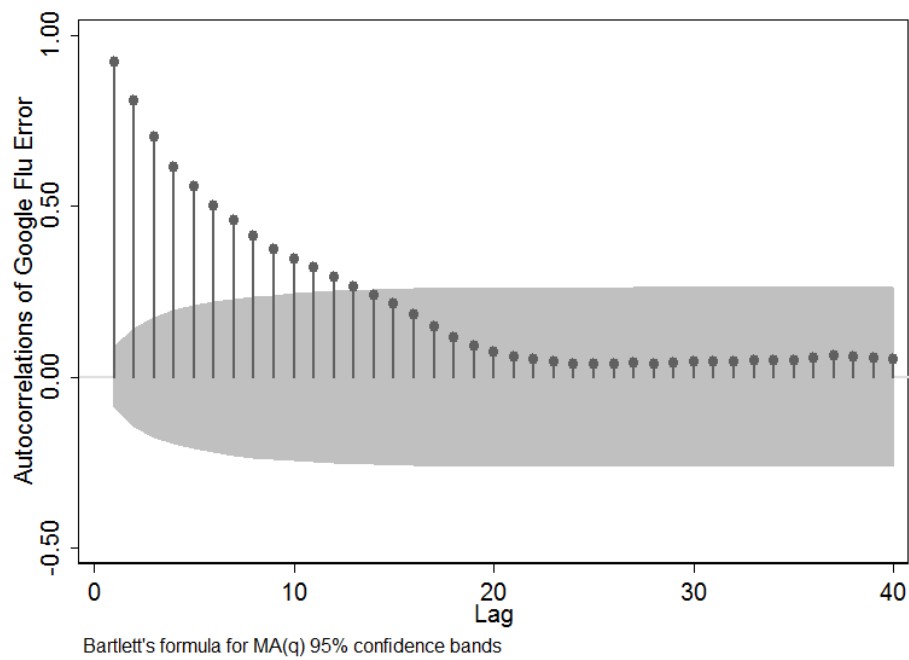


Fig S3: Correlogram of Google Flu Trends Errors

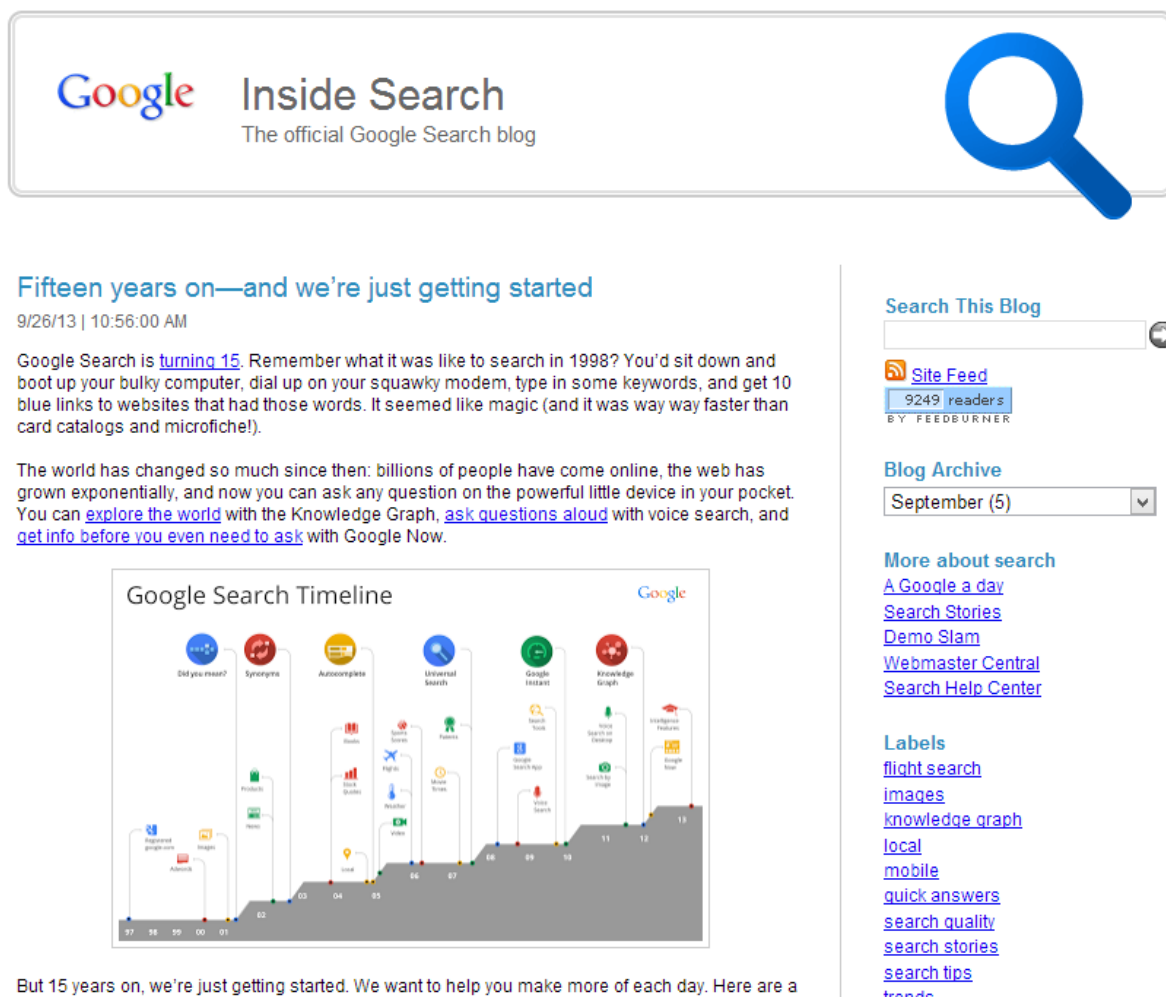


3. Tracking Google's Search Algorithm: An Analysis of Google's Official Search Blog

As we note in the main paper, a full accounting of the evolution of Google's search algorithm is nearly impossible – hundreds of changes are made to the algorithm each year, some major and some minor. We can, however, track major changes using Google's official search blog, which has recorded major events for Google since May 2011 (<http://insidesearch.blogspot.com/>).

The first thing to note from this record is how much Google's search algorithm changes over time. As Fig S4 points out, this is a point of pride for Google, as they note in their most recent blog post, celebrating their 15th anniversary.

Fig S4: Google Blog Discussing How Much Things Have Changed in 15 Years



Source: <http://insidesearch.blogspot.com/2013/09/fifteen-years-onand-were-just-getting.html>

Health related searches have been a particular area of emphasis for Google. Fig S5 shows an announcement from the official Google blog about how the company is consistently running experiments to improve the services provided to their users. In this case, the company was using a short poll to determine the primary reasons why people use different health-related search terms in order to provide users with more relevant search results. We suspect that this determination to provide useful information had the unexpected side-effect of increasing certain searches and throwing off GFT's statistics.

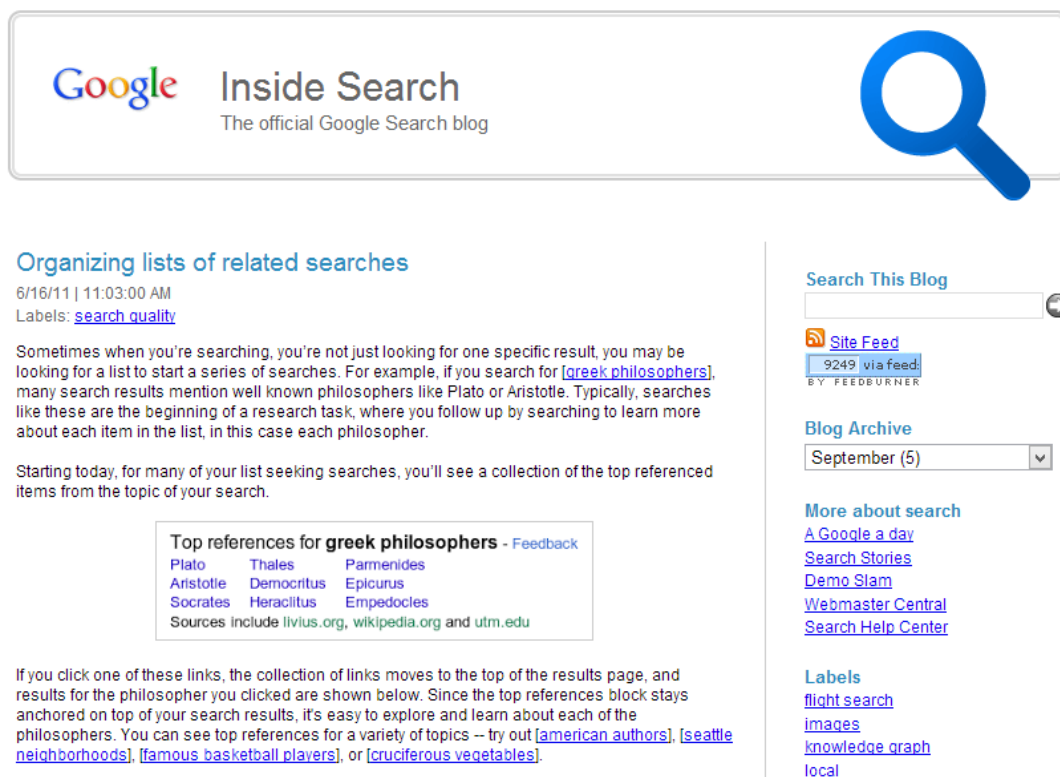
Fig S5: An Example of Google's Consistent Experiments with Health Searches



Source: <http://googleblog.blogspot.com/2009/05/understanding-health-related-searches.html>

An example of the changes made by Google, and also a likely culprit for changing search behavior, is the addition of the related search feature in mid-June 2011, a few weeks before GFT started persistently missing high. This addition provides users with the top referenced searches related to their original search term. For example, a search for “flu” will bring up a number of additional recommended searches about diagnosis and treatment of the flu. As we discuss below, increasing searches for flu treatments are a likely reason for GFT’s growing inaccuracy. Fig S6 records the announcement of this feature.

Fig S6: Google Introduces List of Related Searches Feature



Source: http://insidesearch.blogspot.com/2011/06/organizing-lists-of-related-searches_16.html

Another likely culprit was introduced in February 2012. The health search box was introduced as a way for consumers to search particular symptoms, like “fever” and “cough” and find potential diagnoses of these conditions. Fig S7 records the announcement of this technology and Fig S8 shows the sample search given by Google to demonstrate its use. In Fig S9, we demonstrate this in a more general context of a search for “fever” and “runny nose.” As the reader can see, the top two results are for the flu and the common cold. This seems a likely reason why searches like “cold vs flu” and “cold or flu” seem to spike up in the 2012-2013 flu season.

Fig S7: Google Announcement of Improved Health Searches



Fig S8: Google Sample of New Health Search Box

The image shows a Google search interface. At the top left is the Google logo. To its right is a search bar containing the text "abdominal pain on my right side". Below the search bar, the word "Search" is displayed in red, followed by the text "About 15,300,000 results (0.35 seconds)". On the left side, there is a vertical menu with options: "Everything" (highlighted in red), "Images", "Maps", "Videos", "News", "Shopping", and "More". The main content area displays "Searches related to abdominal pain on my right side". It lists four related terms with brief descriptions: "Appendicitis" (A serious medical condition in which the appendix beco...), "Ovarian cyst" (Any collection of fluid, surrounded by a very thin wall, wi...), "Hernia" (Condition in which part of an organ is displaced and pro...), and "Kidney stone" (A hard mass formed in the kidneys, typically consisting ...). Below these, it lists "Irritable bowel" (A functional bowel disorder characterized by chronic ab...). A note states: "Drawn from at least 10 websites including abdopain.com and wikipedia.org - How this works". Under the "Shopping" category, there is a link titled "Right Side Abdominal Pain ... Causes and Treatment" with the URL "www.abdopain.com/right-side-abdominal-pain.html". Below this link, a snippet of text reads: "Right side abdominal pain is commonly caused by condtions such as ... I have been having pain in my right side lower area for like 3 months now. comes and ...".

Google

abdominal pain on my right side

Search About 15,300,000 results (0.35 seconds)

Everything

Images

Maps

Videos

News

Shopping

More

Searches related to **abdominal pain on my right side**

[Appendicitis](#) A serious medical condition in which the appendix beco...

[Ovarian cyst](#) Any collection of fluid, surrounded by a very thin wall, wi...

[Hernia](#) Condition in which part of an organ is displaced and pro...

[Kidney stone](#) A hard mass formed in the kidneys, typically consisting ...

[Irritable bowel](#) A functional bowel disorder characterized by chronic ab...

Drawn from at least 10 websites including [abdopain.com](#) and [wikipedia.org](#) - [How this works](#)

[Right Side Abdominal Pain ... Causes and Treatment](#)
[www.abdopain.com/right-side-abdominal-pain.html](#)

Right side abdominal pain is commonly caused by condtions such as ... I have been having pain in **my right side** lower area for like 3 months now. comes and ...

Source: <http://insidesearch.blogspot.com/2012/02/improving-health-searches-because-your.html>

Fig. S9: Google Search for “runny nose” and “fever”

The screenshot shows a Google search interface with the query "runny nose and fever" in the search bar. The search results page displays "About 3,530,000 results (0.28 seconds)". An advertisement for Triaminic is shown, with a red arrow pointing to the "Searches related to runny nose and fever" section. This section lists related terms: Flu, Common cold, Bacterial infection, Sinusitis, and Pneumonia, each with a brief description. Below this, there are two links to Yahoo Health symptom search pages for "Runny nose, Fever, Stuffy nose" and "Runny nose, Cough, Fever".

Google

runny nose and fever

Web Images Maps Shopping More Search tools

About 3,530,000 results (0.28 seconds)

Ad related to **runny nose and fever**

Runny Nose Symptoms
www.triaminic.com/
Triaminic® Products Provide Relief To Your Child's Runny Nose

Searches related to runny nose and fever

Flu	An infectious disease caused by rna viruses of the family...
Common cold	Viral infectious disease of the upper respiratory system, ...
Bacterial infection	The detrimental colonization of a host organism by a fore...
Sinusitis	Inflammation of a nasal sinus
Pneumonia	Lung inflammation caused by bacterial or viral infection, i...

Drawn from at least 10 websites, including webmd.com and wikipedia.org - [How this works](#)

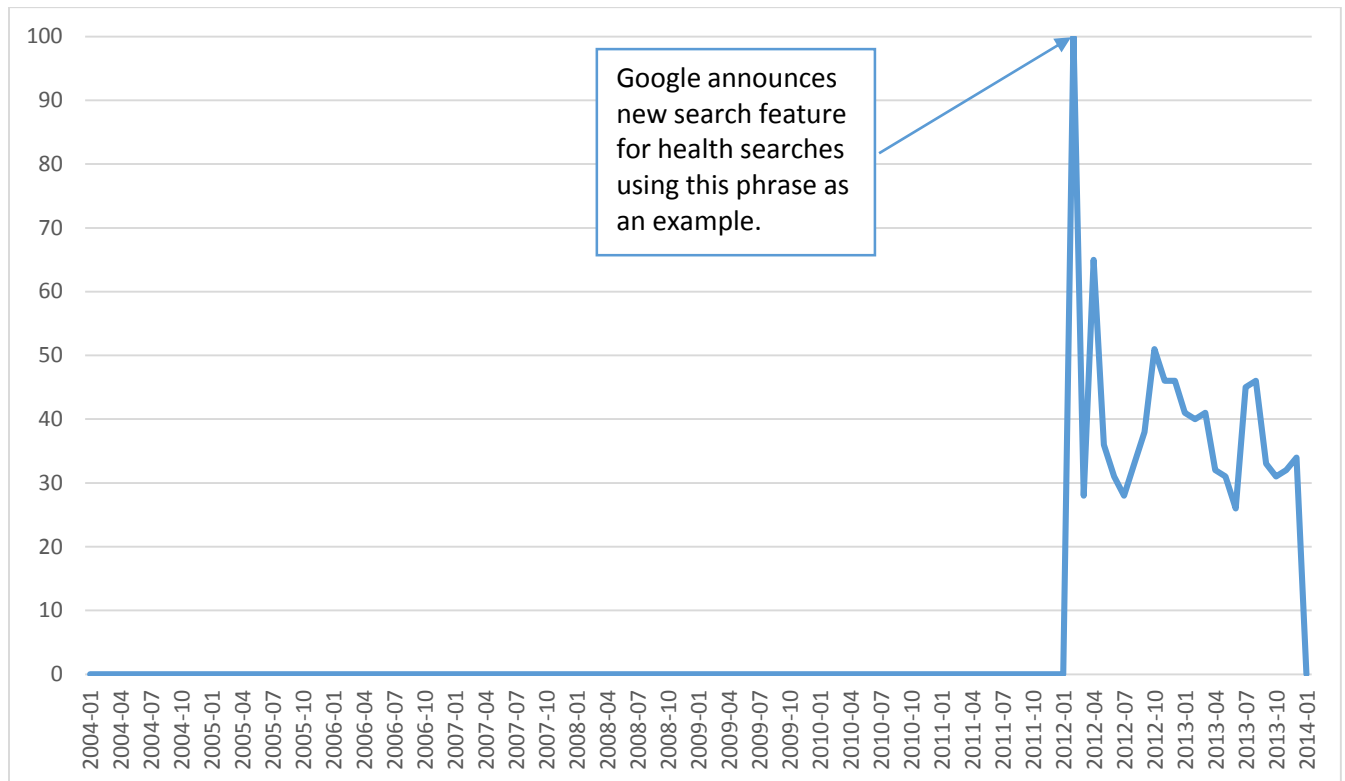
[Symptom Search for Runny nose, Fever, Stuffy nose on Yahoo Health](http://health.yahoo.net/symptomsearch?addterm=Runny+Nose...Fever...Nose)
health.yahoo.net/symptomsearch?addterm=Runny+Nose...Fever...Nose
Discover 9 possible causes for Runny nose, Fever, Stuffy nose including Common Cold Overview, Flu Overview and Sinus Infections (Sinusitis)

[Symptom Search for Runny nose, Cough, Fever on Yahoo Health](http://health.yahoo.net/symptomsearch?addterm=Runny+Nose...Cough...Fever)
health.yahoo.net/symptomsearch?addterm=Runny+Nose...Cough...Fever
Discover 11 possible causes for Runny nose, Cough, Fever including Common Cold Overview, Flu Overview and Sinus Infections (Sinusitis)

Note: Search for “runny nose and fever” (see search line) conducted from Boston, MA in October 2013. Results last replicated on February 26, 2014. Searches were conducted using both Google Chrome and Mozilla Firefox browsers.

This particular announcement also had an added effect – it demonstrates the ability for Google itself to influence the relative prevalence of search term use. As one might suspect, searches for “abdominal pain in my right side” are relatively rare. The relative prevalence of this search, however, sparks markedly right after the announcement using this term as an example. This spike is documented in Fig S10.

Fig S10: Searches for “abdominal pain on my right side”

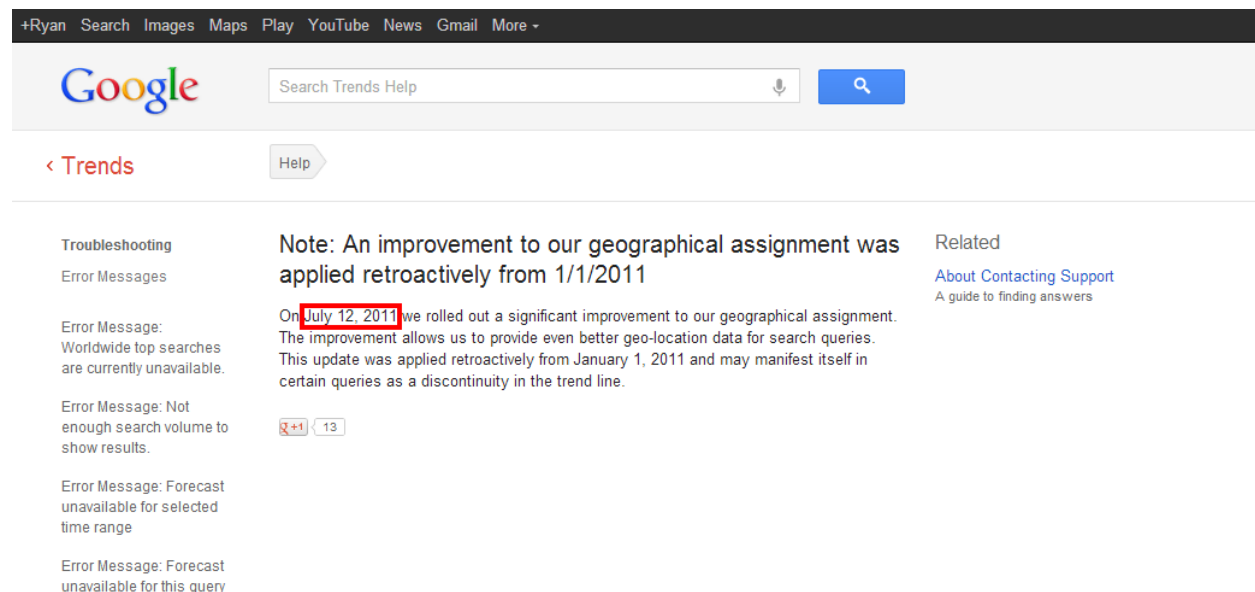


Source: Google Trends (<http://www.google.com/trends/>). Downloaded data available in replication materials.

Note: Google Trends results are based on the proportion of searches for the phrase relative to total searches. The reported results are scaled so 100 represents the highest relative volume of searches for the phrase. All others are reported as the percentage relative to the highest search volume. The peak, in this case, is more than four standard deviations above the average for this search term. The peak is reached in February 2012. Google announced the new search feature on February 13, 2012 (see above).

These are far from the only changes that might have affected GFT's performance. As noted above, the performance of GFT after 2011 was markedly different between regions. This is important, as the model was originally developed based on regional flu time series, rather than the national time series (Ginsberg et al. 2009, p. 2). Interestingly, it is in July 2011 that Google announced a substantial improvement in its geolocation abilities for search that were retroactively applied to all searches after January 1, 2011. This change was not documented on the official search blog, but is noted in certain searches in Google Trends. We document this in Fig S11.

Fig. S11: Documentation of Change in Geographic Resolution in 2011



Note: As indicated by the red highlight, this change was originally implemented on July 12, 2011. This note was originally posted at

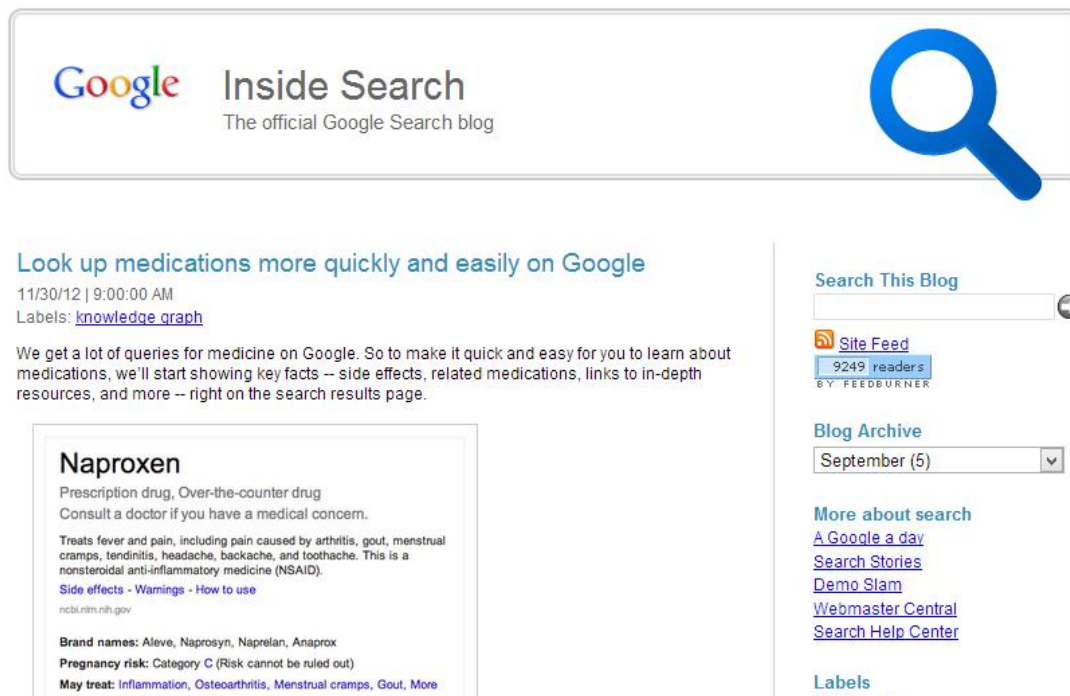
<https://support.google.com/trends/?hl=en&rd=2#topic=4365599>, but has subsequently been removed. Documentation of this change can be found in:

L. Rossignol, C. Pelat, B. Lambert, A. Flahault, E. Chartier-Kastler, and T. Hanslik. A Method to Assess seasonality of Urinary Tract Infections Based on Medication Sales and Google Trends. *PLOS One*, doi: 10.1371/journal.pone.0076020 (2013).

And also: <http://www.epiphanysearch.co.uk/blog/2011/07/warning-google-makes-insights-for-search-data-useless/>

Even smaller changes can also affect how GFT will work in the future. At the end of 2012, Google announced new features to help people with finding information on different medications, their uses and potential interactions. From a business and customer service perspective, the goal is to increase the usage of Google for information on medications, but given that “robitussin” is listed as one of the example search terms for GFT by Google, and there are likely other treatments present, an increase in the relative frequency of searches for medication might influence GFT results. We document this change in Fig S12.

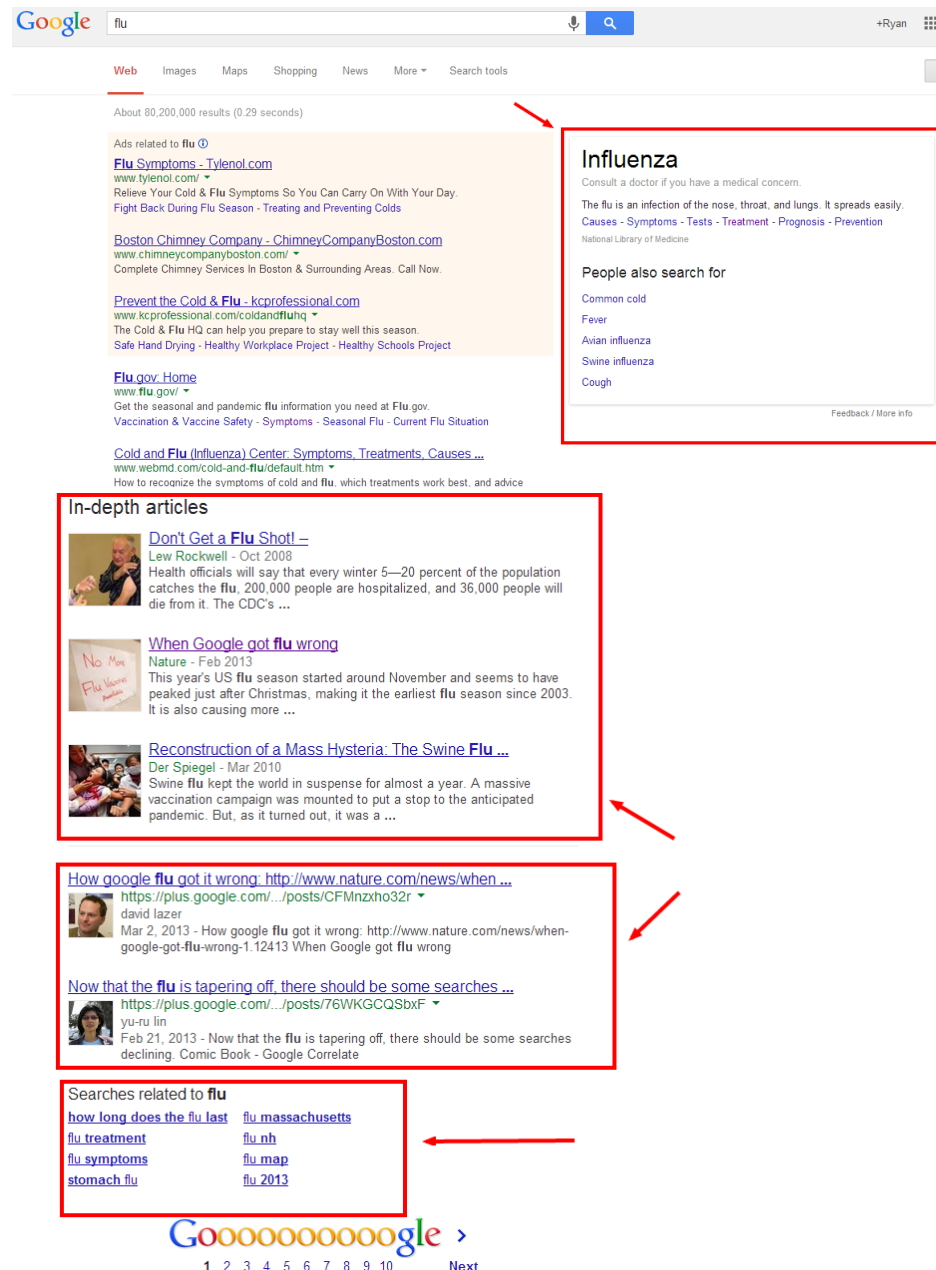
Fig. S12: Google Feature to Look Up Medication Information



Source: <http://insidesearch.blogspot.com/2012/11/look-up-medications-more-quickly-and.html>

The last three figures present some summary issues related to Google search. Fig S13 provides a visual summary of many of the changes made in the last few years to how Google presents results. As the reader can see, a standard search for “flu” produces a number of new choices that may affect consumer behavior. Among these are a very noticeable side bar on influenza that includes common searches as well as treatment and symptom information, a set of in-depth articles of varying quality, links to posts by Google+ contacts, and the related searches option discussed above.

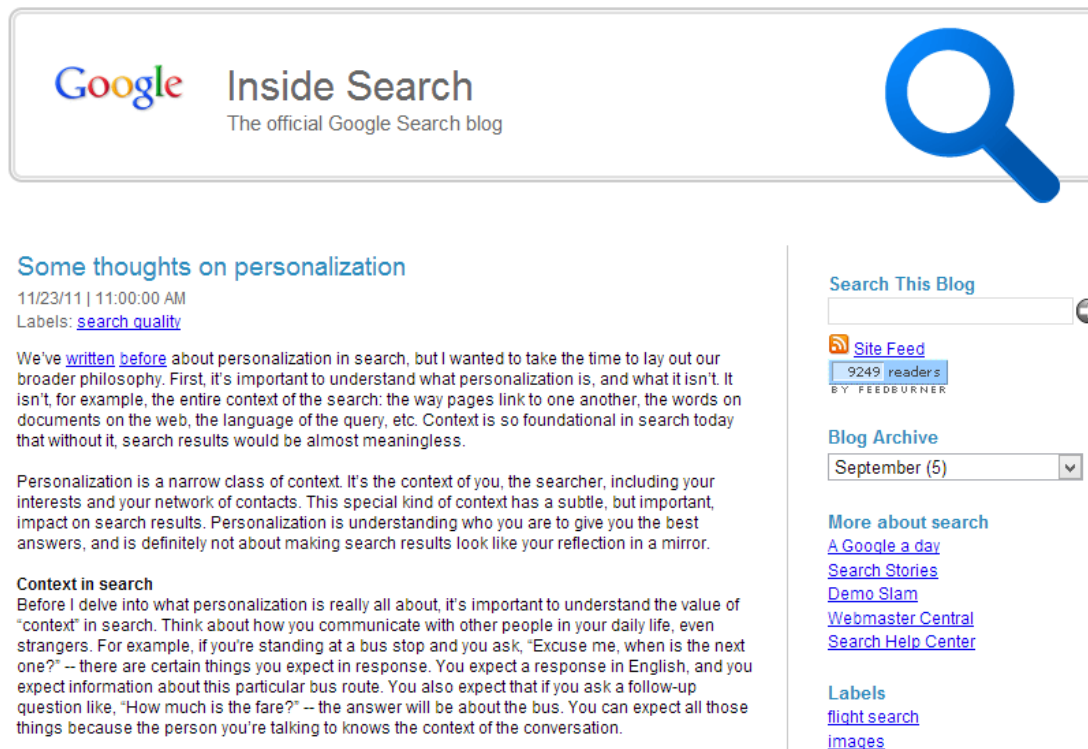
Fig. S13: Examples of Other Google Search Changes for “flu”



Note: Search for “flu” (see above search line) conducted from Boston, MA in October 2013. Results last replicated on February 26, 2014. Searches were conducted using both Google Chrome and Mozilla Firefox browsers.

Figure S14 deals with search personalization, a topic for which the academic community needs a much greater understanding. This is an excerpt from a blog post about Google's personalization. It discusses how they attempt to use context in order to improve searches. Such context is likely to change the results and the types of searches recommended for consumers.

Fig. S14: Google Discusses Personalization



Source: <http://insidesearch.blogspot.com/2011/11/some-thoughts-on-personalization.html>

Finally, Fig S15 is included to point out that Google may not only experience “blue team” dynamics in its search patterns. Indeed, some companies hire engineers specifically to reconstruct Google’s search algorithm so that they can be at the top of search results. There are now companies who, as part of their marketing services will try to change placement on search results, and, of course, Google also sells this space with its sponsored content pattern. The item shown in this illustration is the introduction of hot searches. Much like Twitter’s trending hashtags, this would seem to be an area ripe for manipulation by campaigns and companies interested in getting media attention for their brand.

Fig. S15: Google Announcement of Search Trends Feature

New ways to explore what's trending on Google

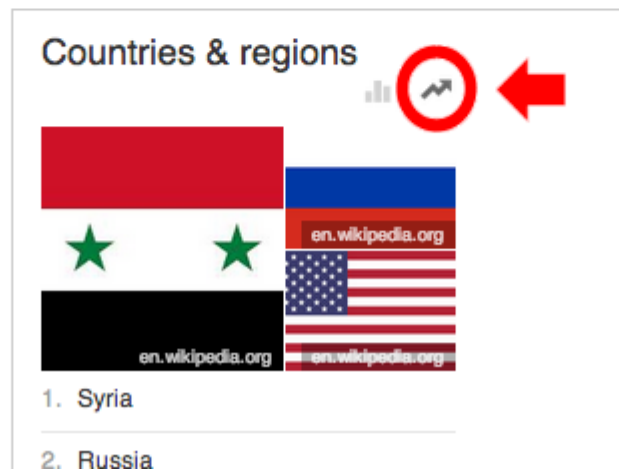
9/12/13 | 8:55:00 AM

Whether you're looking for trending celebrities, a monthly recap of what's hot, or power tools to make your own discoveries about what's piquing the world's curiosity - today you'll find new features in Google Trends to make it easier to explore hot topics in Google Search.

Trending Top Charts.

In May we [added a new feature](#) to Google Trends called “Top Charts,” where you can explore real-world people, places and things ranked by overall search interest in the United States (with more countries coming soon). These “Most Searched” lists span dozens of areas from athletes to cities to cocktails. We've heard great feedback from people who want “Trending” lists -- not just what's most searched overall, but what's spiking compared with usual search volumes. Starting today, you can explore these new Trending Top Charts for a number lists across entertainment, sports, politics and more.

For example, while it may come as no surprise that the United States is the most searched country among people in the U.S., it's more interesting that Syria and Russia were the two top trending countries last month. To see the new “Trending” charts, click the arrow icon at the top of any supported [Top Chart](#).



Source: <http://insidesearch.blogspot.com/2013/09/new-ways-to-explore-whats-trending-on.html>

4. Search Dynamics – Which Search Terms Likely Led to GFT’s Inaccuracy and the Problem of Replicating GFT’s Results.

We start with the last topic listed in the title – the difficulty replicating GFT’s results. We attempted in vain to find the 45 search terms that were utilized by GFT. They were not listed in the supplemental materials for any of the published articles on GFT, nor were they available from any online sources that we could locate. The original *Nature* article listed 12 categories of search terms, but only listed examples of search terms that were not utilized (e.g. high school basketball). A later article, published by *PLOS One*, did give some examples of the search terms utilized [10]. The eight examples listed were: symptoms of bronchitis, pneumonia (sic), fever, early signs of the flu, robitussin, influenza a, amoxicillin, and strep throat (p. 2). We plotted the time trend for these terms in Fig S16 using results from Google Trends. It is notable that, with the exception of “influenza a,” none of these examples correlated well with either the GFT or the CDC data. While it is possible that this is a result of the aggregation mechanism used for GFT, the fact that the technology designed to allow others to access the power that made GFT produces results so at odds with the information provided by the authors gave us pause. We cannot be sure if these are really examples used to produce GFT or if the examples were purposefully misleading. A later article discussing GFT’s 2013 update seems to suggest that they were indeed purposefully misleading [11].

Fig. S15: Google Trend Searches for Terms Indicated in PLOS One Article



Note: Data acquired by typing search terms into Google Trends (<http://www.google.com/trends/>) and downloading the associated data. The downloaded data are available in the SOM/SOM4/SFig15 folder of the replication materials.

Also troubling is that none of those sample terms, at least in their verbatim form, came up when we used Google Correlate to find the most related terms with different time series. Tab S3 lists the top correlated terms found by Google Correlate for the GFT data, the GFT data post-2011, the CDC data, and the CDC data pre-2009 (what is received when one clicks the link “...match the pattern of actual flu activity (this is how we made Google Flu Trends!)”). As one can note from this table, none of the eight example terms are listed in the top-50 for any of these time series. The table also lists the relative ranks of the different terms and a rough categorization that we put together for analysis. This provides the first clues for which terms were responsible for GFT’s errors.

Tab S3: List of Search Terms and Classification

Search Text	Rank for GFT	Rank for CDC	Rank for CDC Pre-2009	Rank for GFT Post-2011	Classification of Search
influenza type a		1	1		Term for Influenza
flu duration		2	3	37	General Information
flu fever		3	5	39	Flu Symptoms
treating flu		4	23		Remedies/Treatments
braun thermoscan		5	40		Flu Diagnosis
fever flu		6	33		Flu Symptoms
flu recovery		7	10	34	General Information
flu vs. cold		8	24		Flu Diagnosis
cold or flu	2	9	11	18	Flu Diagnosis
treating the flu		10			Remedies/Treatments
Oscillococcinum	29	11	34		Related Diseases
flu versus cold		12	49		Flu Diagnosis
flu remedies	6	13	35	23	Remedies/Treatments
cold versus flu		14			Flu Diagnosis
human temperature		15	26		Flu Diagnosis
contagious flu	17	16	31		General Information
type a influenza		17	46		Term for Influenza
flu or cold	3	18	13	48	Flu Diagnosis
flu contagious	1	19	4	8	General Information
Thermoscan		20			Flu Diagnosis
flu cough		21			Flu Symptoms
influenza incubation period		22			General Information
duration of flu	40	23			General Information
cold vs flu	8	24	39	24	Flu Diagnosis
influenza a		25			Term for Influenza
low body temperature		26			Flu Symptoms
flu headache		27		42	Flu Symptoms

Search Text	Rank for GFT	Rank for CDC	Rank for CDC Pre- 2009	Rank for GFT Post- 2011	Classification of Search
flu complications		28			Complications
flu stomach		29			Flu Symptoms
cold and flu symptoms		30			Flu Symptoms
flu and fever		31		31	Flu Symptoms
cold vs. flu		32			Flu Diagnosis
treatment for flu	9	33	25		Remedies/Treatments
treatment of flu		34			Remedies/Treatments
ear thermometer		35			Flu Diagnosis
how long does the flu last?		36			General Information
flu in children		37			General Information
influenza incubation		38			General Information
flu length		39			General Information
length of flu		40			General Information
type a flu		41			Term for Influenza
getting over the flu		42		29	Remedies/Treatments
treat flu	7	43	16		Remedies/Treatments
robitussin ac		44			Remedies/Treatments
treatment for the flu	48	45			Remedies/Treatments
influenza symptoms		46	38		General Information
what is influenza		47			General Information
flu care		48			Remedies/Treatments
Expectorant		49			Remedies/Treatments
cold symptoms		50			Related Diseases
flu germs	33				Term for Influenza
cure flu	27				Remedies/Treatments
how to treat flu	18				Remedies/Treatments
how to get rid of the flu	49			1	Remedies/Treatments
get rid of the flu	36		45	2	Remedies/Treatments
fever reducer			22		Remedies/Treatments
i have the flu	37		47	26	General Information
flu treatment			20		Remedies/Treatments
dangerous fever			27		Flu Symptoms
remedies for the flu	12			7	Remedies/Treatments
medicine for the flu	25			5	Remedies/Treatments
how long does it take to get the flu				33	General Information
how to cure the flu	38			19	Remedies/Treatments
flu in toddlers				27	General Information

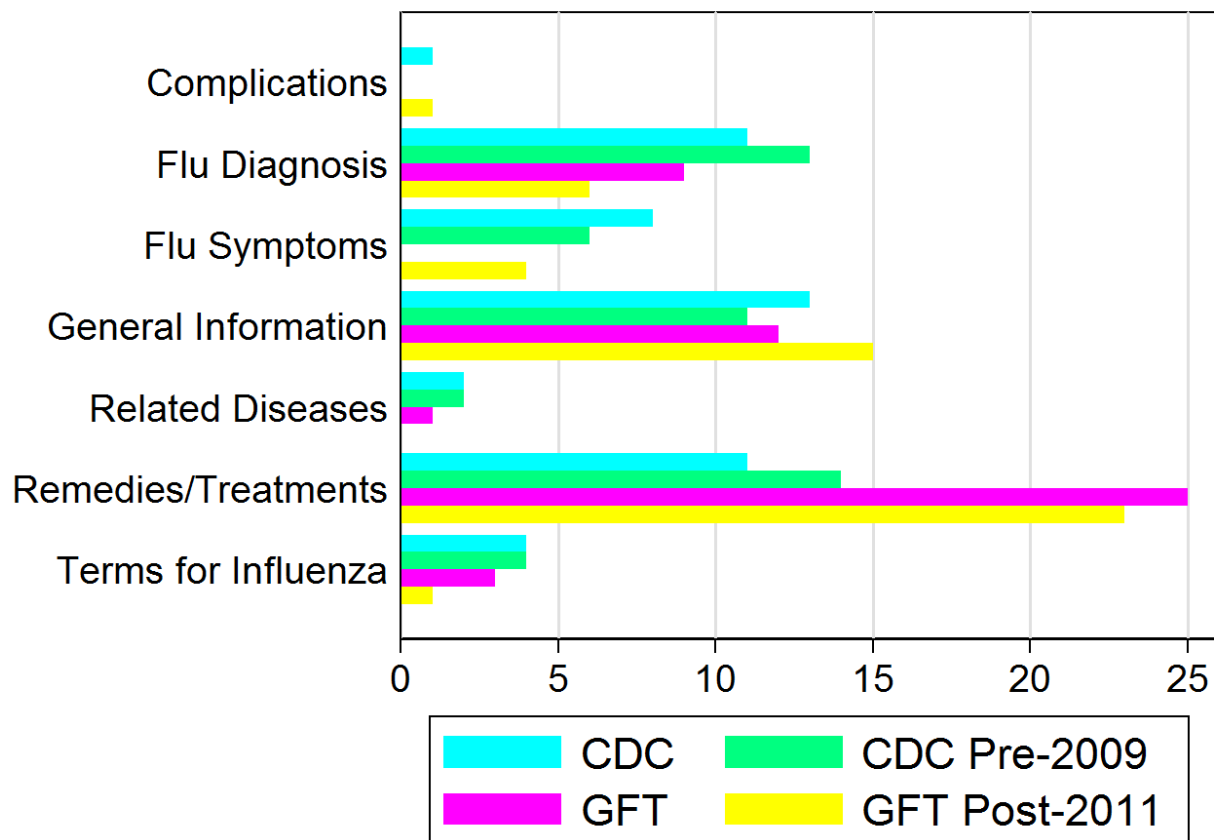
Search Text	Rank for GFT	Rank for CDC	Rank for CDC Pre-2009	Rank for GFT Post-2011	Classification of Search
symptoms of flu	13		2		General Information
treat the flu	4		6	36	Remedies/Treatments
Bronchitis			50		Related Diseases
reduce fever			19		Remedies/Treatments
when is the flu contagious	32			28	General Information
flu a				35	Term for Influenza
the flu	20		28		Term for Influenza
home remedies for flu				14	Remedies/Treatments
flu medicine	19		12	3	Remedies/Treatments
influenza a and b	23		30		Term for Influenza
medicine for flu	34				Remedies/Treatments
taking temperature			48		Flu Diagnosis
cure the flu	16			13	Remedies/Treatments
signs of the flu			8		Flu Diagnosis
flu and cold	10				Flu Diagnosis
do i have the flu	28				Flu Diagnosis
flu treatments			37		Remedies/Treatments
how long is the flu contagious	46				General Information
is flu contagious	30		15	38	General Information
how long does the flu last			32		General Information
fight the flu	44				Remedies/Treatments
normal body			14		Flu Diagnosis
home remedies for the flu	39			11	Remedies/Treatments
best flu medicine				20	Remedies/Treatments
flu relief				44	Remedies/Treatments
signs of flu			42		Flu Diagnosis
symptoms of influenza	41				Flu Diagnosis
how to treat the flu	5		7	40	Remedies/Treatments
flu swab				10	Flu Diagnosis
how long is flu contagious			36	47	General Information
what to do for the flu	42				Remedies/Treatments
how long are you contagious				9	General Information
how long am i contagious				15	General Information
body temperature			17		Flu Symptoms
flu home remedies	24			21	Remedies/Treatments
normal body temperature			44		Flu Symptoms
get rid of flu				4	Remedies/Treatments

Search Text	Rank for GFT	Rank for CDC	Rank for CDC Pre-2009	Rank for GFT Post-2011	Classification of Search
over the counter flu medicine	31			50	Remedies/Treatments
the flu virus	43				General Information
how to get rid of flu				16	Remedies/Treatments
is the flu contagious	35		18	32	Remedies/Treatments
how long do flu symptoms last				41	General Information
cold and flu	15			43	Flu Diagnosis
pregnant with the flu				22	Complications
incubation period for the flu	47				General Information
exposure to flu	50				General Information
how long does flu last	21		43		General Information
flu vs cold	11		21		Flu Diagnosis
natural flu remedies	45				Remedies/Treatments
symptoms of the flu			9		Flu Diagnosis
am i contagious				6	General Information
how long flu last				45	General Information
difference between cold and flu	22				Flu Diagnosis
flu how long				25	General Information
flu test				49	Flu Diagnosis
how to get over the flu	26			17	Remedies/Treatments
best medicine for flu				46	Remedies/Treatments
remedies for flu	14		29	12	Remedies/Treatments
fever cough			41		Flu Symptoms
flu without fever				30	Flu Symptoms

Note: Search terms compiled from Google Correlate (<https://www.google.com/trends/correlate>) using national-level time series data. The “replication” of GFT using Google Correlate (CDC pre-2009) is located at (<https://www.google.com/trends/correlate/search?e=id:20xKcnNqHrk&t=weekly#>). Data used for this table is available in the replication materials in folder SOM/SOM4/STab3.

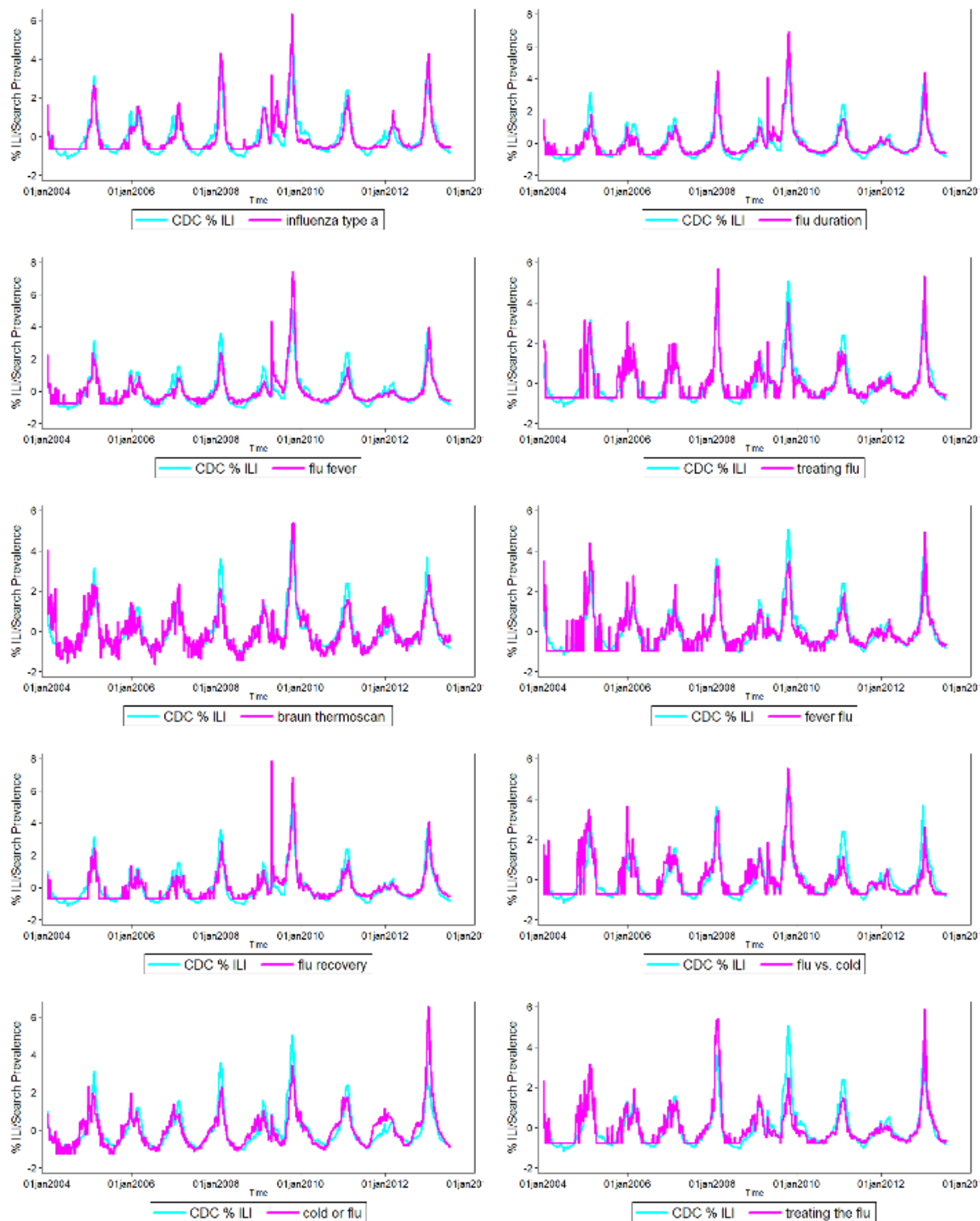
One notable pattern that is immediately apparent is that the number and rank of terms dealing with “treatment for flu” and its variants is much higher for the GFT data than it is for the CDC data. Indeed, Fig S16 shows that the number of terms in the category of “Remedies/Treatments” that fall within the top-50 for GFT time series is about double the number of treatment related terms that fall in the top-50 for the CDC time series. This would seem to suggest that an increased prevalence of searches for remedies and treatments was a likely culprit in throwing GFT’s count off.

Fig. S16: Category of Top 50 Search Terms From Google Correlate



Further evidence for this hypothesis can be found by looking at the pattern for the top-10 terms correlated with the CDC's full time series of data (Fig S17). Most of the terms that fall into the top-10 continue to resemble the CDC's data throughout the time period, including after 2011, when GFT started to persistently miss high. Two types of terms are the exception to this rule. The first are terms related to remedies and treatments for the flu. Searches for "treating flu" and "treating the flu" are both very high for the 2012-2013 flu season. The other search result that is noticeably high for this period is the search for "cold or flu." Looking back on the table above, there are a large number of variations of "cold or flu," "flu or cold," etc. that seem to show up with the GFT data and not with the CDC data. After searching through a number of additional terms, an increase in these two search terms seems the most likely explanation for GFT's increasing errors. As we note in the main text of the article, trying to pin down the exact causes is impossible without at least having the original search terms – and this is why we emphasize replication – but these two seem very likely candidates for explaining what happened.

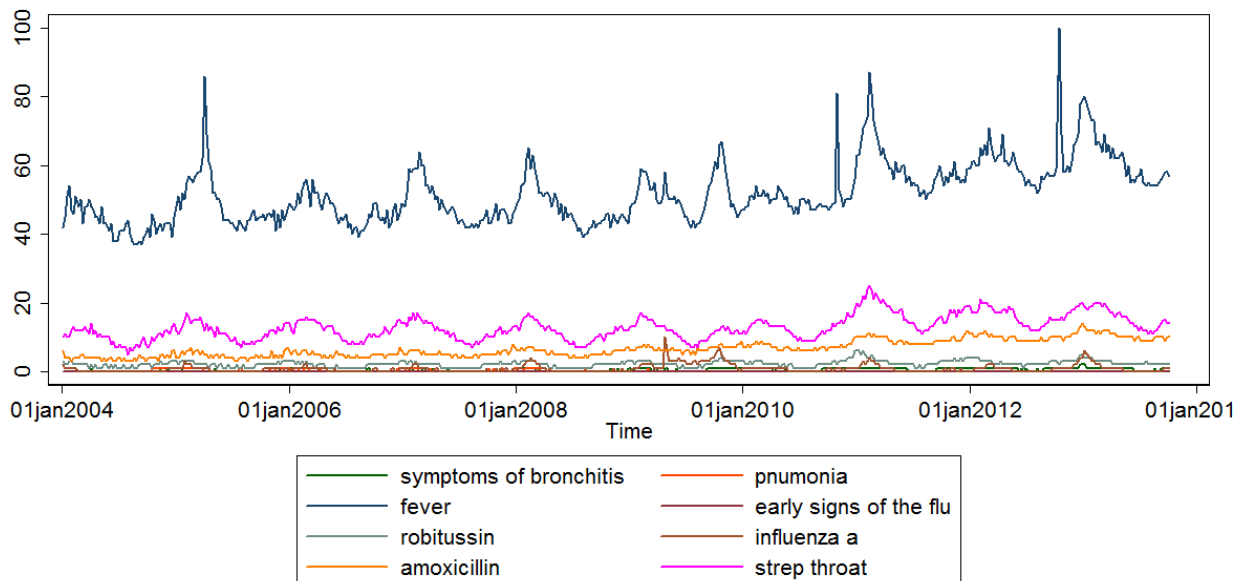
Fig S17: Top-10 Google Correlate Search Terms with CDC Data



Note: Data acquired from the best correlated search terms with CDC percent ILI data as reported by Google Correlate (<https://www.google.com/trends/correlate>). Replication data available in SOM/SOM4/SFig17 folder of replication materials.

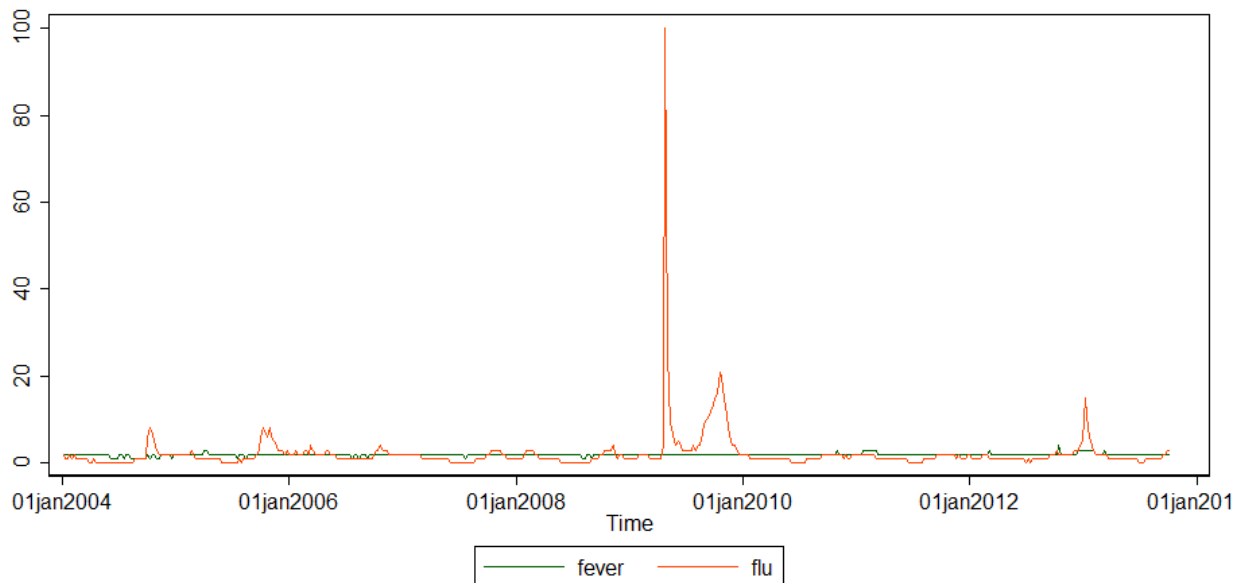
Another factor that may have contributed to the error is that the search terms used may not be very prevalent in the first place. Again, taking the authors at their word, we explore the relative prevalence of the eight example search terms from the *PLOS One* article. Using “fever,” by far the most prominent term, as the baseline, we can see in SFig 18 that most of these terms are searched relatively rarely. What this means is that, while they may have had a high correlation, they are unlikely to be robust to changes in search patterns. A relatively small change in their frequency may have a large effect on their relative frequency. SFig 19 shows that even the most searched example term, “fever,” is still a relatively rare search compared to general searches about “flu.” The general search for “flu” should also give us pause when talking about the common explanation for GFT’s miss – that there was massive media attention given to the flu in 2012-2013, resulting in more flu-related searches. First, searches for “flu” peaked in 2009, which is what one might expect given that this is when the swine flu panic took place. We should also note that not all flu-related searches spiked in 2012-2013, only a select few did. And these select few do not seem to be a random sample of the overall search terms. This does not seem to fit well with the media panic story.

Fig S18: Relative Frequency of Search Terms from GFT’s PLOS One Article



Note: Data derived from search prevalence of terms as reported by Google Trends (<http://www.google.com/trends/>). Data in SOM/SOM4/SFig15 folder of replication materials.

Fig. S19: Relative Frequency of Searches for “Fever” Versus “Flu”



Note: Data derived from search prevalence of “fever” and “flu” as reported by Google Trends (<http://www.google.com/trends/>). Data available in the SOM/SOM4/SFig19 folder of the replication materials.

5. A Note on GFT's 2008 Model

In the main text, we made reference to GFT's 2008 algorithm, which was substantially modified after the 2009 influenza a H1N1 pandemic. The data from this algorithm is no longer available from Google Flu. A previous study made use of data downloaded in August 2009 to conduct an analysis of the algorithm's performance and documents its problems during the 2009 pandemic. Today, however, the only available data is from March 2009 as made available through the internet archive (aka "The Wayback Machine"). It can be downloaded from the following link: <http://web.archive.org/web/20090303211715/http://www.google.org/about/flutrends/download.html>.

We make this data available, along with the global estimates, in the replication materials, so they can be preserved. They can be found in folder /SOM/SOM5/.

6. Lag Models Using CDC Data 3 and 4 Weeks Out

One criticism that could be leveled against both our models using lagged CDC data and those that have been analyzed previously is that there is unmodeled data vintaging within the CDC data (i.e. the data are subsequently corrected by the reporting agency and the originally reported data is usually not made available). Similar vintaging is common in economic measures, which are regularly retrospectively revised.

If early errors in CDC data reporting are randomly distributed, then this would not pose much of a problem for our results. Nonetheless, it is worthwhile for us to conduct robustness checks using longer lags to make sure we are not receiving artificially high results due to using already corrected data.

Unsurprisingly, the best-fitting models with only 3+ week lags are a little different from those in the main paper. For this section, our lagged CDC variable model is specified as

$$flu_t = \alpha + \beta_1 flu_{t-3} + \beta_2 flu_{t-4} + \beta_3 flu_{t-5} + \sum_{i=1}^{52} \gamma_i week_{it}$$

where flu is the CDC estimated percent of doctors' visits for ILI and $week$ is coded 1 if the week is number n in a year.

The combination of GFT and the CDC data is specified as

$$flu_t = \alpha + \beta_1 gflu_t + \beta_2 flu_{t-3} + \beta_3 flu_{t-4} + \beta_4 (gflu_{t-3} - flu_{t-3}) + \sum_{i=1}^{52} \gamma_i week_{it}$$

where $gflu$ is the GFT estimate of ILI.

Again, the regression models were initially specified using only data prior to the September 2009 launch of the new GFT and were tested on data produced subsequently. The out-of-sample predictions are made on a rolling basis, with predictions for the next time step ($t+1$) based on estimates from data on all previous time periods ($t_i \leq t$).

Despite using further out lags, these models still outperform the GFT data by itself. For the out-of-sample time period, mean absolute error for the GFT data, by itself, is 0.486. For the model using only lagged CDC data, the mean absolute error is 0.412. For the model using a combination of GFT's estimates and the lagged CDC data, the mean absolute error is 0.303, or about a 38% decrease in error.

Code for replicating these results can be found in the SOM/SOM6 folder of the replication materials. The code is in the file SOM6 (Replication Code).do. The data is located in the Manuscript folder and is ParableOfGFT(Replication).dta. All calculations were done in Stata.

7. Full Results and Sensitivity Tests

For space reasons, full results of the models using lagged CDC data and combinations of the CDC data and GFT data were not included. In this section, we present not only the models underlying the graphs in the main text, but also conduct a general sensitivity test to show how sensitive, if at all, these results are to different specifications. The general sensitivity test should reassure readers about potential over-fitting due to the inclusion of higher-order autoregressive variables or about difficulties in comparing models with different numbers of features.

In the table below, we compare the following model specifications:

$$flu_t = gflu_t \quad (1)$$

$$flu_t = \alpha + \beta_1 flu_{t-2} \quad (2)$$

$$flu_t = \alpha + \beta_1 gflu_t + \beta_2 (gflu_{t-2} - flu_{t-2}) \quad (3)$$

$$flu_t = \alpha + \beta_1 flu_{t-2} + \sum_{i=1}^{52} \gamma_i week_{it} \quad (4)$$

$$flu_t = \alpha + \beta_1 flu_{t-2} + \beta_2 flu_{t-3} + \sum_{i=1}^{52} \gamma_i week_{it} \quad (5)$$

$$flu_t = \alpha + \beta_1 flu_{t-2} + \beta_2 flu_{t-3} + \beta_3 flu_{t-4} + \sum_{i=1}^{52} \gamma_i week_{it} \quad (6)$$

$$flu_t = \alpha + \beta_1 gflu_t + \beta_2 (gflu_{t-2} - flu_{t-2}) + \sum_{i=1}^{52} \gamma_i week_{it} \quad (7)$$

$$flu_t = \alpha + \beta_1 gflu_t + \beta_2 flu_{t-2} + \beta_3 (gflu_{t-2} - flu_{t-2}) + \sum_{i=1}^{52} \gamma_i week_{it} \quad (8)$$

$$flu_t = \alpha + \beta_1 gflu_t + \beta_2 flu_{t-2} + \beta_3 (gflu_{t-2} - flu_{t-2}) + \beta_4 (gflu_{t-3} - flu_{t-3}) + \sum_{i=1}^{52} \gamma_i week_{it} \quad (9)$$

where flu is the CDC estimate of ILI and $gflu$ is the GFT estimate of ILI, $week$ is a binary variable indicating week of observation, and β and γ are estimated regression coefficients. Specifications 1, 6 and 9 are the ones presented in Fig 2 of the main manuscript.

Data and code for running all of these models and calculating their mean absolute error (MAE) and root mean squared error (RMSE) can be found in the /Manuscript/ folder of the replication materials and can be produced using the ParableOfGFT(Replication).dta data with the Parable of Google Flu (Replication Code).do code file. The code for calculating statistical significance from this data can be found in the /SOM/SOM7/ folder of the replication materials (SOM7(Replication Code).do). All calculations were done in Stata.

As was done in the main paper, all estimates are originally trained on the data prior to the release of the 2009 GFT update (September 6, 2009). From that point on, predictions are done on a rolling basis, similar to how such a system would likely be deployed in the real world. Predictions in time step $t+1$ are based on estimates from data on all previous time periods, $t_i \leq t$.

As the reader can see clearly in the table below, all of the specifications are an improvement on the GFT estimate in the out-of-sample (after September 6, 2009) period. The difference between all other models (2-9) and the GFT model are also statistically significant.

The reader can also clearly see that the models which combine the information from GFT with the information from the CDC (3, 7, 8, and 9) perform much better than their counterparts that utilize either the CDC data or the GFT data alone. While the models presented in the main paper are over-specified, they are not overfit, as their performance in this out-of-sample data is still an

improvement over less highly specified models using the same data, although these differences are not always statistically significant.

Measures of statistical significance will vary based on the loss function for the errors. In the main paper, we report mean absolute error (MAE), which is resistant to outliers. Tab S4 reports these differences along with a sign test for statistical significance. The sign test evaluates the hypothesis that the median difference between the out-of-sample errors is zero. While not a very powerful test, it makes very few assumptions about the distribution of the data and is resistant to outliers. Other scholars have suggested that this makes it an attractive significance test for comparing forecast errors across non-nested models (e.g. Diebold and Mariano 1995, p. 254-255). The test is also somewhat conservative in this context, as the results cannot be driven by a few large GFT misses during the last flu season.

A more common approach in regression forecasts is to report the root mean squared error (RMSE) and assume a squared error loss function (i.e. larger misses are more harmful than smaller misses). This is somewhat consistent with epidemiologists' concern with estimating the size of the peak of flu season, which is also where forecast errors are highest. In this situation, a rather intuitive, but somewhat cumbersome, test of the null hypothesis that the mean squared forecast errors of the two models are equal (an F-test) is suggested by Ashley et al. (1980; see also Ashley 1981; Ashley 2003). These results, along with the RMSE are reported in Tab S5. This test's utility is somewhat limited by its assumption of squared error loss and it being only asymptotically justified (while the sign test is an exact finite-sample test), but the size of our out-of-sample cases is well above what Ashley (2003) finds to be appropriate.

Ashley, R., C.W.J. Granger, and R. Schmalensee. 1980. "Advertising and Aggregate Consumption: An Analysis of Causality." *Econometrica* 59: 817-858.

Ashley, R. 1981. "Inflation and the Distribution of Price Changes Across Markets." *Economic Inquiry* XIX: 650-660.

Ashley, R. 2003. "Statistically Significant Forecasting Improvements: How Much Out-Of-Sample Data is Likely Necessary?" *International Journal of Forecasting* 19: 229-239.

Diebold, Francis X. and Roberto S. Mariano. 1995. "Comparing Predictive Accuracy." *Journal of Business and Economic Statistics* 13: 253-265.

Tab S4: Comparison of Models – MAE with Sign Test for Significance

	1	2	3	4	5	6	7	8	9
1	--	0.148 (0.072)	0.237 (0.000)	0.134 (0.000)	0.169 (0.001)	0.174 (0.000)	0.243 (0.000)	0.251 (0.000)	0.254 (0.000)
2		--	0.089 (0.000)	-0.014 (0.128)	0.021 (0.038)	0.027 (0.053)	0.096 (0.000)	0.103 (0.001)	0.106 (0.000)
3			--	-0.103 (0.019)	-0.068 (0.000)	-0.063 (0.000)	0.006 (0.782)	0.014 (0.489)	0.016 (0.580)
4				--	0.035 (0.782)	0.041 (1.000)	0.109 (0.013)	0.117 (0.027)	0.120 (0.027)
5					--	0.005 (0.000)	0.074 (0.000)	0.082 (0.000)	0.084 (0.000)
6						--	0.069 (0.000)	0.077 (0.000)	0.079 (0.000)
7							--	0.008 (0.213)	0.010 (0.489)
8								--	0.003 (0.072)
9									--

Note: Values in the cells are the difference in mean absolute error (MAE) between models with p-values (2-tailed) based on a sign test in parentheses. Error is estimated for the out-of-sample period (post September 6, 2009) using the dynamic procedure outlined above.

Tab S5: Comparison of Models – RMSE with F-Test for Significance

	1	2	3	4	5	6	7	8	9
1	--	0.371 (0.000)	0.435 (0.000)	0.335 (0.000)	0.410 (0.000)	0.413 (0.000)	0.449 (0.000)	0.482 (0.000)	0.487 (0.000)
2		--	0.064 (0.072)	-0.036 (0.139)	0.039 (0.137)	0.041 (0.110)	0.078 (0.032)	0.111 (0.002)	0.116 (0.001)
3			--	-0.101 (0.007)	-0.026 (0.418)	-0.023 (0.463)	0.014 (0.261)	0.046 (0.000)	0.051 (0.000)
4				--	0.075 (0.001)	0.078 (0.001)	0.115 (0.003)	0.147 (0.000)	0.152 (0.000)
5					--	0.003 (0.416)	0.039 (0.239)	0.072 (0.023)	0.077 (0.014)
6						--	0.037 (0.269)	0.069 (0.026)	0.074 (0.016)
7							--	0.033 (0.000)	0.038 (0.000)
8								--	0.005 (0.176)
9									--

Note: Values in the cells are the difference in root mean squared error (RMSE) between models with p-values (2-tailed) based on an F-test in parentheses. Error is estimated for the out-of-sample period (post September 6, 2009) using the dynamic procedure outlined above.

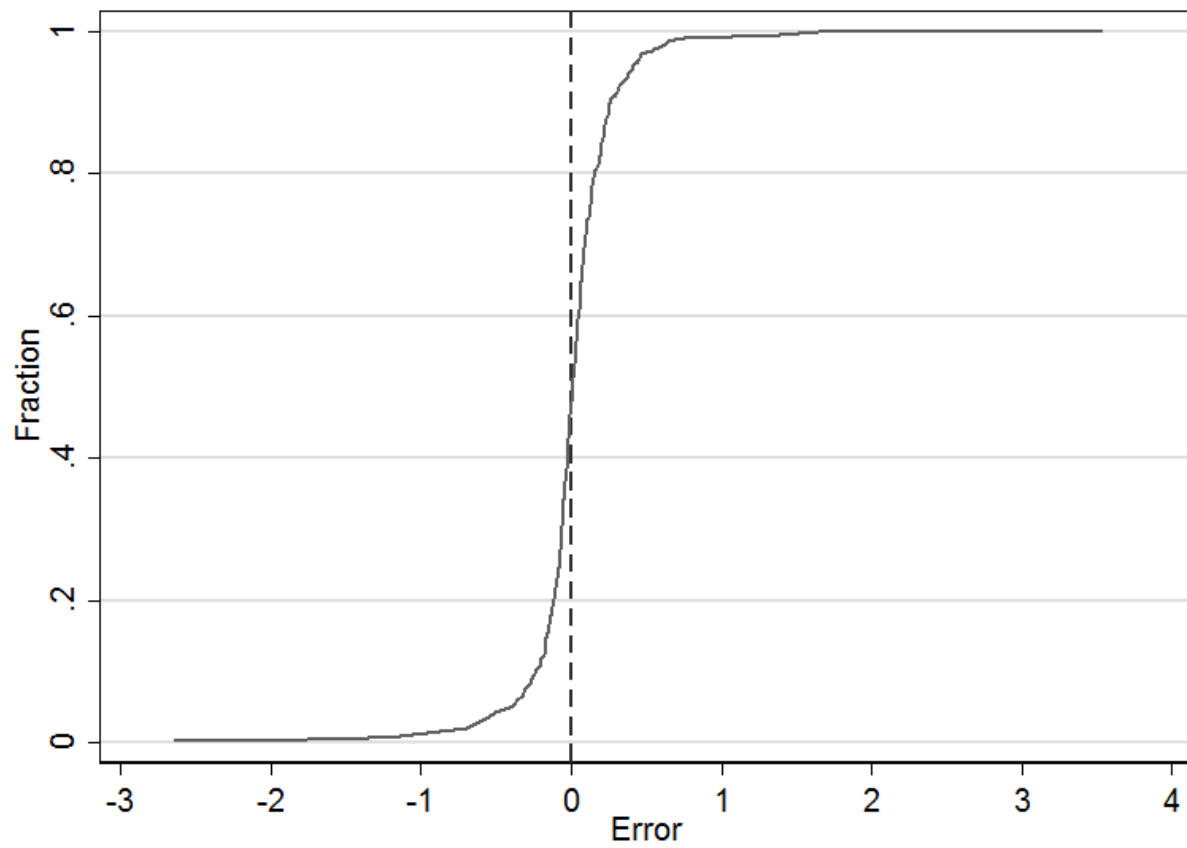
8. Comparison With Nonlinear Models

In this section of the SOM, we utilize a more flexible, nonlinear, technique for modeling the CDC ILI estimate. Linear models have dominated this particular area of analysis (see e.g. Ginsburg et al. 2009; Cook et al. 2011; Goel et al. 2010), and have clear advantages in simplicity and interpretability. Linear models, however, may be insufficient in the presence of even simple nonlinearities (Weigend and Gershenfeld 1993: 16). This section serves as a check on whether we can gain substantial ground through use of simple nonlinear models.

Replication for all of these results can be done using the code in the /SOM/SOM8/ folder of the replication materials. The code SOM8(Replication Code).do will allow the user to replicate the error plots in conjunction with the ParableOfGFT(Replication).dta data file in the /Manuscript/ folder using Stata. The neural network models were run using the “nnet” package in the R statistical programming environment. This code is also available in the /SOM/SOM8/ folder and is labelled SOMpt8(Replication Code).r.

We begin by conducting the sanity checks suggested by (Weigend and Gershenfeld 1993: 44; see also Smith, same volume: 311 and 336). The first check is to look at the distribution of errors to see if they are uniform or if there are a few large and unusual outliers. Figure Fig S26 shows a cumulative distribution plot by the size of errors from the combination of GFT and CDC data presented in the main paper in the out-of-sample data. While this seems to show the pattern Weigend and Gershenfeld are concerned about (very small errors along with a few large outliers), these outliers are very rare.

Fig S26: Cumulative Distribution of Errors from Combined GFT and CDC Model



Similarly, Fig S27 and S28 show the errors plotted against the actual CDC ILI data and the predicted values from the model respectively. Here again, the errors look relatively uniform with a few relatively large (though not as large as for GFT by itself) exceptions.

Fig S27: Plot of Errors from Combined GFT and CDC Model Against Actual CDC ILI Values

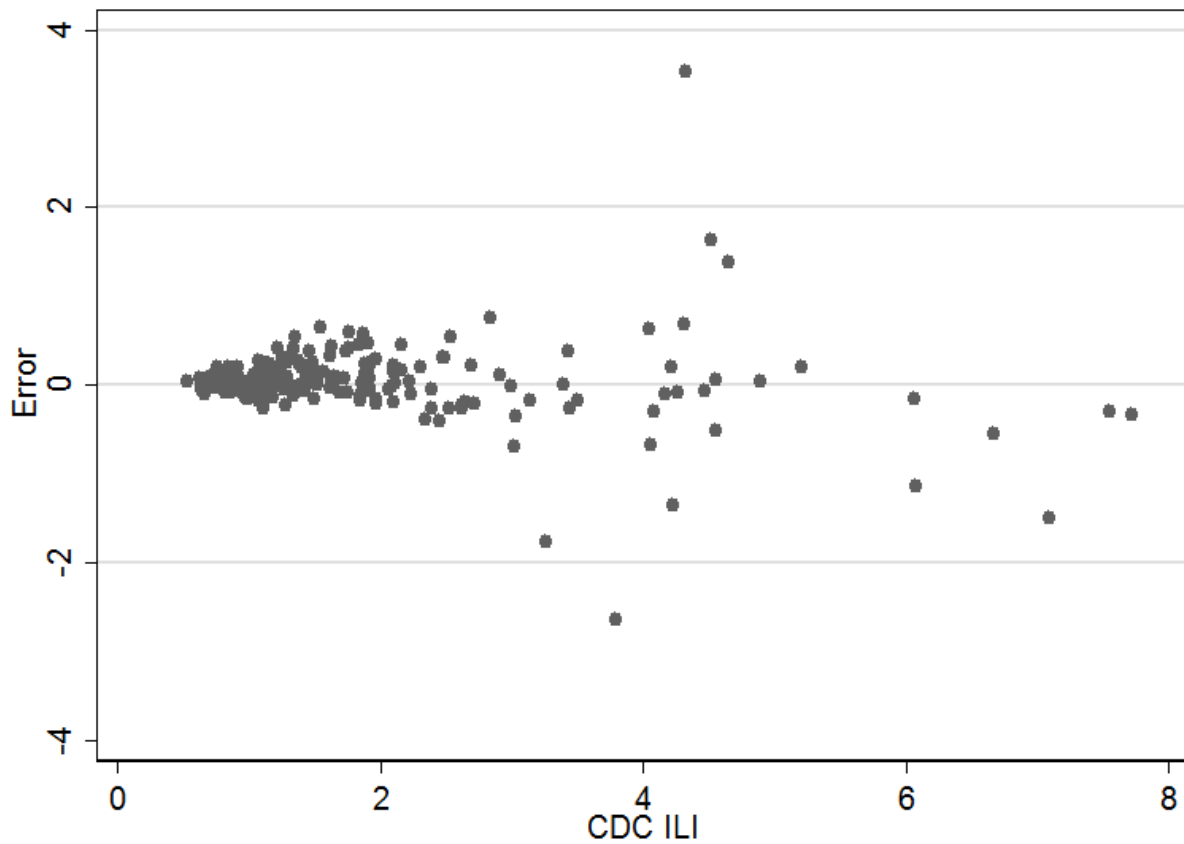
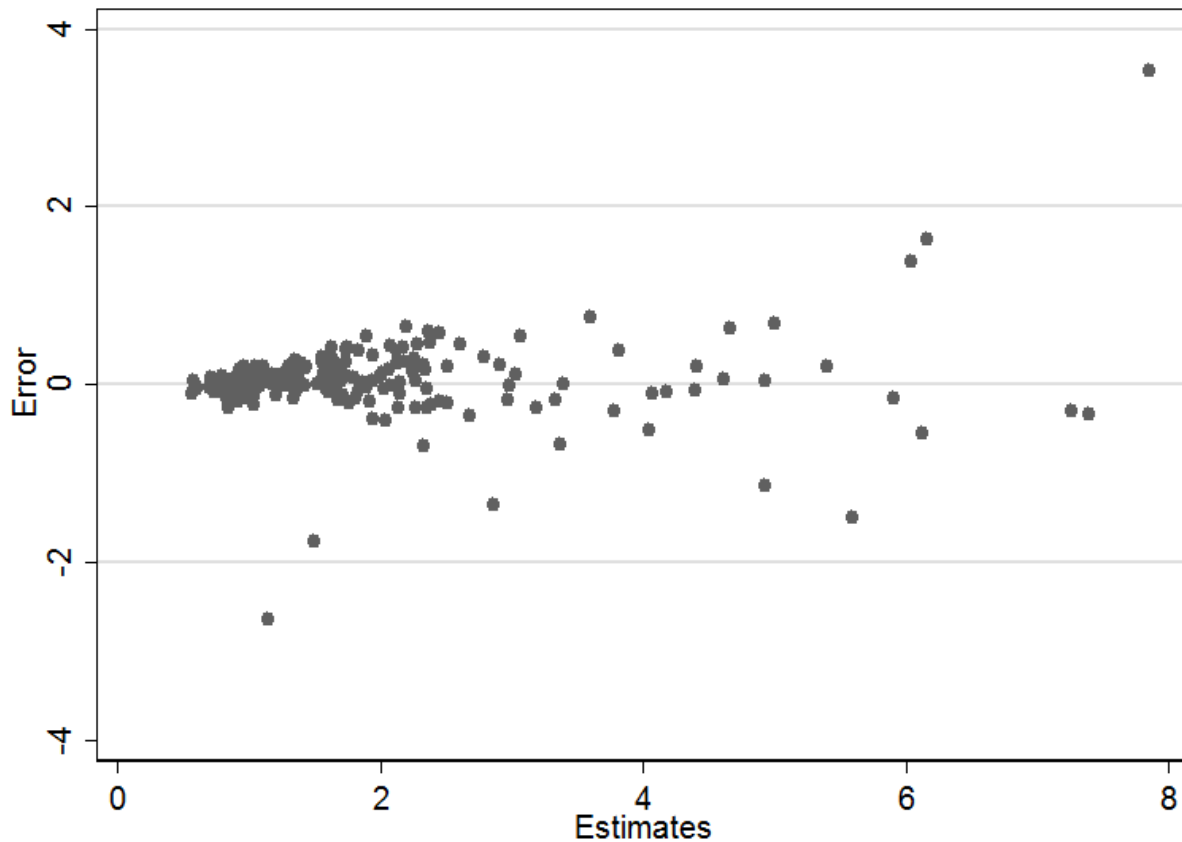


Fig S28: Plot of Errors from Combined GFT and CDC Model Against Predicted Values



These figures would seem to indicate that there is potential nonlinearity in the system, but the sparsity of these outliers and their location in the out-of-sample data mean that a nonlinear estimator may not produce any significant improvement.

To test this we utilize a simple feed-forward neural network model with a single hidden layer architecture. This choice of model architecture is only one of many possible options. It does, however, provide a good first cut estimation of the gains to be made from relaxing the linearity assumptions of the models in the main paper.

The greater flexibility of the neural network model requires us to specify an additional tuning parameter, namely the number of nodes in the hidden layer. To set this parameter in a manner that avoids overfitting, we utilize leave-one-out cross-validation on the training data. We test across a range of possible sizes for the hidden layer. Preference is given to the simplest architecture that cannot be significantly improved with an additional node.

Another issue with neural network models (as is common in more flexible models) is that initial values, which are drawn using a quasi-random number generator, can subtly influence the results, especially when the amount of data available for training is not especially large. We thus

evaluate all models over 50 runs to ensure that our results are not due to the behavior of random number selection.¹

As in the main paper, the in-sample period is prior to the fielding of the 2009 GFT model, with the out-of-sample period estimated on a rolling basis – mimicking the deployment of GFT.

The inputs for our three models are:

$$\{gflu_t, flu_{t-2}, (gflu_{t-2} - flu_{t-2})\} \quad (1)$$

$$\{gflu_t, flu_{t-2}, (gflu_{t-2} - flu_{t-2}), week_{it}\} \quad (2)$$

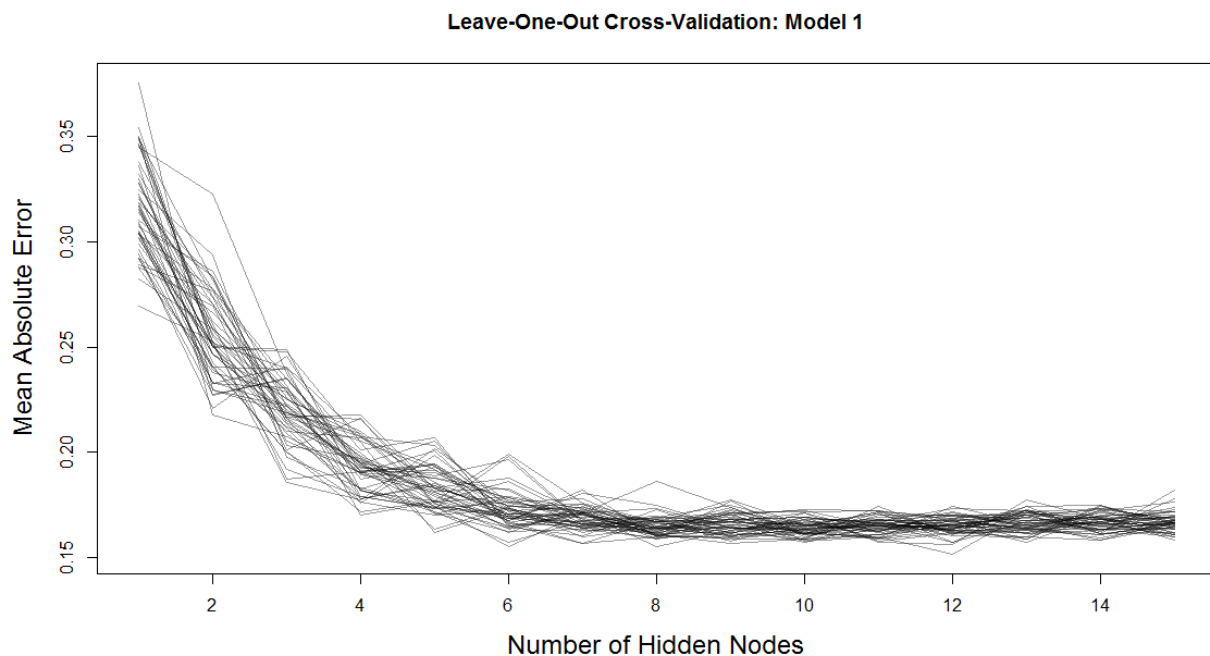
$$\{gflu_t, flu_{t-2}, (gflu_{t-2} - flu_{t-2}), (gflu_{t-3} - flu_{t-3}), week_{it}\} \quad (3)$$

where *flu* is the CDC estimate of ILI and *gflu* is the GFT estimate of ILI.

The variables used in models 2 and 3 are the same as those used in model 8 and 9 in the linear models. The variables for model 1 are a subset of those in model 2 without the seasonality component.

The results of the leave-one-out cross-validation for model 1 are reported in Fig S29. The lines chart the mean absolute error (MAE) across each run. Error seems to stabilize at its lowest level with 10 hidden nodes.

Fig S29: Leave-One-Out Cross-Validation for Model 1

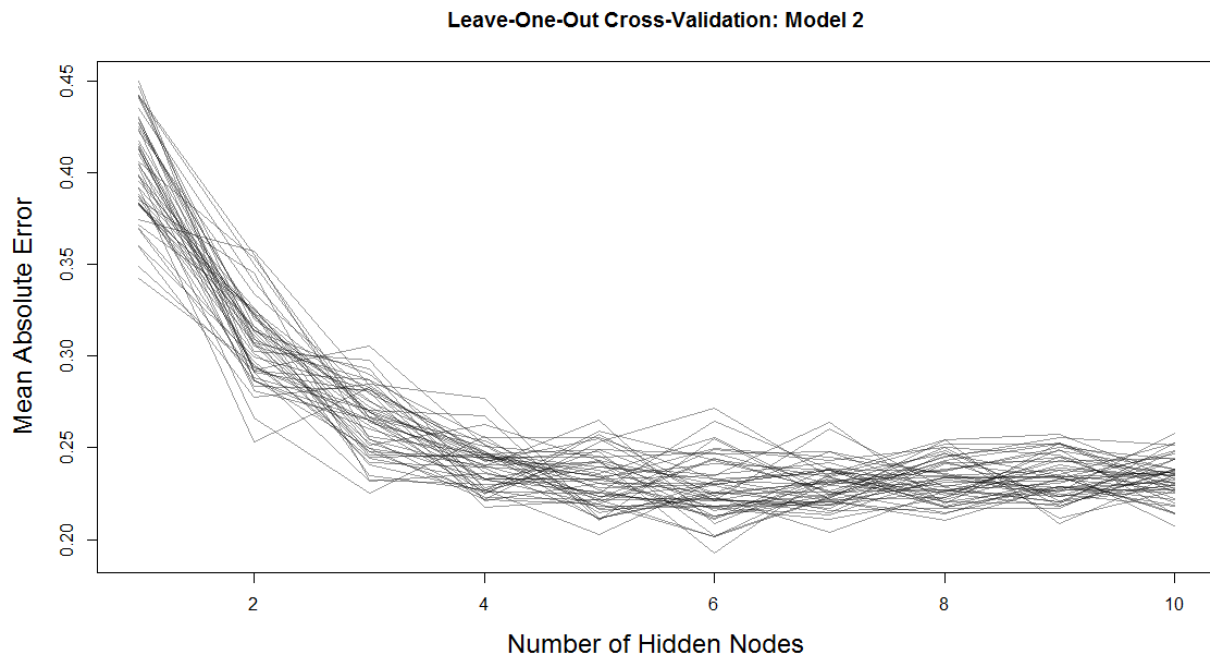


¹ In some applications this is handled by setting a seed for the random number generator. We eschew this approach here because it is important to know about the stability of the model's predictions in order to avoid a favorable/unfavorable seed influencing the interpretation of the results.

After 50 runs of the neural network model with 10 nodes in the hidden layer, we evaluate the mean absolute error of the model. MAE for this neural networks model is 0.249, with a range between 0.215 and 0.318. This is not a substantial improvement over the linear models presented in the main paper.

Model 2 adds in the seasonality component. Here again, we start with the leave-one-out cross-validation to establish the number of hidden nodes. Fig S30 seems to suggest stabilization around 5 nodes in the hidden layer.

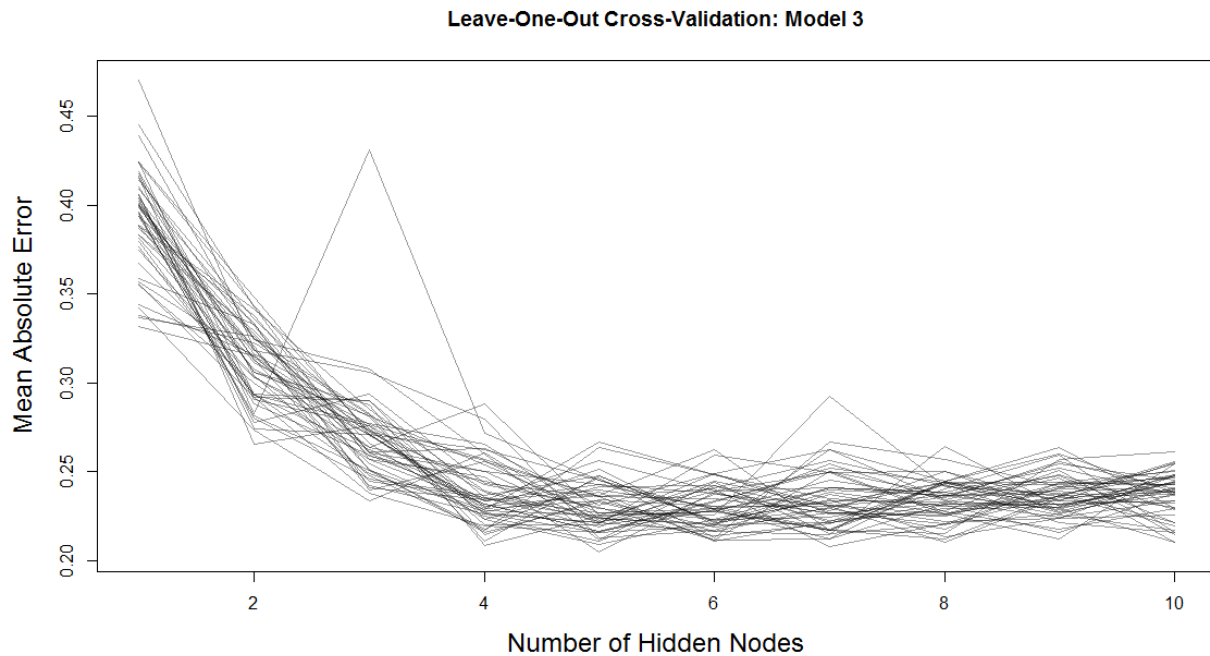
Fig S30: Leave-One-Out Cross-Validation for Model 2



Here again, however, this model does not, on average, outperform the simple linear model in the main paper. The MAE for this model is 0.240 with a range from 0.217 to 0.297. It is clearly an improvement on the model without seasonality, but it is still not an improvement on the linear model.

The final model uses all of the variables from the combined CDC and GFT model in the main paper. Fig S31 shows the results of the cross-validation procedure. As with model 2, it appears to stabilize around five hidden nodes.

Fig S31: Leave-One-Out Cross-Validation for Model 3



While this model comes close to replicating the error rate of the linear model, it is very difficult to distinguish from the linear model (the average difference is measured in 1/1000ths of a point). The average MAE across 50 runs of the model is 0.236, with a range from 0.198 to 0.273.

Clearly, the results presented here cannot reject the possibility that a nonlinear model will perform better (i.e. perhaps a more complex neural network architecture would improve the results), but it does suggest that the linear models presented in the main paper are very difficult to outperform through simple relaxation of the linearity assumption. We should also note that linear models likely perform well in our situation because we are only predicting a couple of steps into the future (technically, we are “nowcasting” – trying to figure out current levels when measurement lags the event). In studies that have attempt to predict influenza levels months in advance, like those cited in the main paper, nonlinearity looms much larger. It is also possible that, since there were only a few large outliers that occurred primarily in the out-of-sample data, nonlinear prediction will become more pertinent over time.