



EUROPEAN CENTRAL BANK

EUROSYSTEM

A SPOT in the dark: Using AI to assess financial stability risks

11th BoF and ESRB Joint Conference on AI
and Systemic Risk Analytics

04/06/2026

***Domenic Kellner, Jan Hannes Lang,**
Marek Rusnák (ECB) and Lukas Nagy (University of Marburg)*



Disclaimer

The views expressed in this paper are those of the authors and do not necessarily reflect those of the European Central Bank

Outline

- 1 Introduction
- 2 Method and Data
- 3 Results
- 4 Conclusion

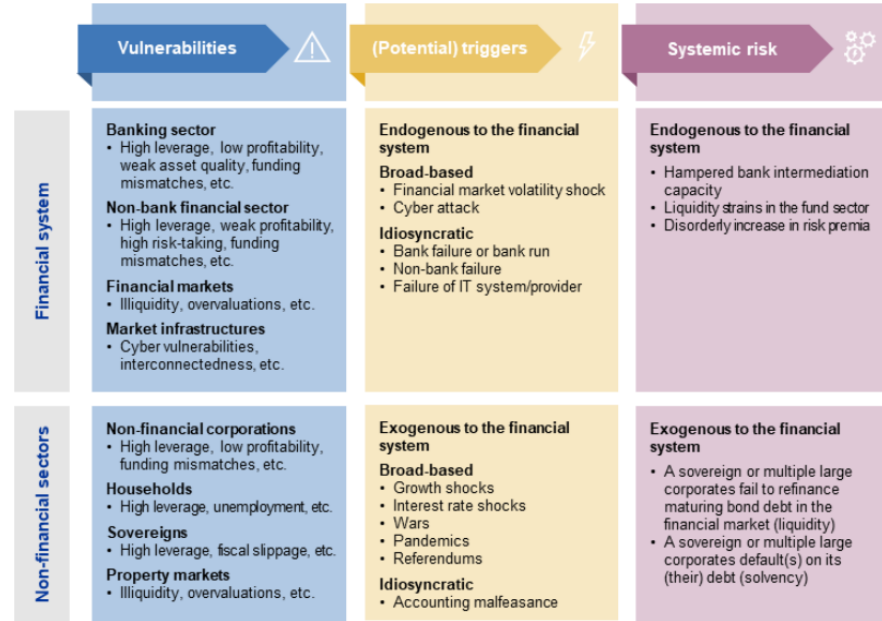
1

Introduction

Motivation: Financial stability monitoring of triggers mainly qualitative

- Financial stability risks have two components:
 - Vulnerabilities:** leverage, overvaluation, ...
 - Triggers:** bank failures, wars, pandemics, ...
- Monitoring of vulnerabilities is well developed
- Monitoring of triggers remains largely qualitative**
- Our motivation:** systematically **assess potential trigger events** using news and **AI**

ECB Financial Stability Review framework: vulnerabilities, triggers and risks



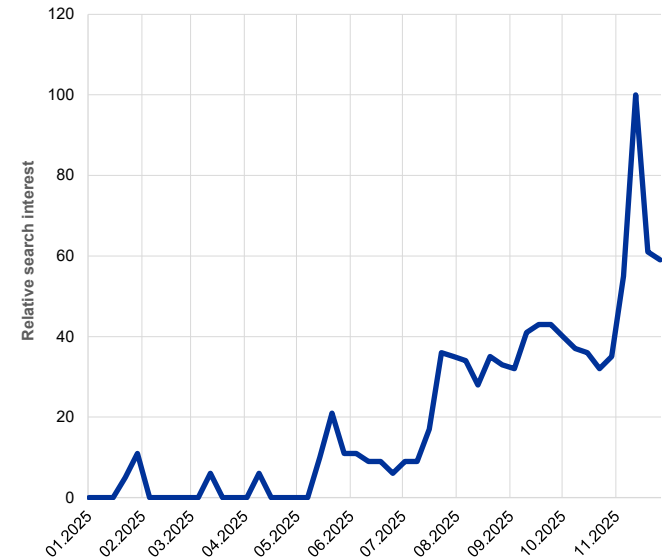
Source: Figure A.2 from Nov 2024 ECB Financial Stability Review Special Feature by Fell et al. "[Communication for financial crisis prevention: a tale of two decades](#)"

Why AI and text-based methods for financial stability analysis?

- Financial stability (FS) risks often materialise following rare but severe trigger events, making early detection challenging
- Many FS indicators reflect vulnerabilities (e.g. credit, RRE prices, balance sheets) and are released with time lags
- News and narratives contain information about potential trigger events but have been hard to monitor systematically so far
- AI enables scalable and systematic analysis of unstructured text, supporting timely identification of financial stability triggers

Increasing interest in the topic of Financial Stability and AI

(Relative search interest)



Sources: Google trends data

Notes: Search term is Financial Stability and AI, worldwide. Numbers represent search interest relative to the highest point on the chart for the given region and time. A value of 100 is the peak popularity for the term. A value of 50 means that the term is half as popular. A score of 0 means there was not enough data for this term. This is relative search interest rather than absolute search volume.

From dictionary indices to LLM-based indicators

| | Traditional Natural Language Processing (NLP) | Transformer- and LLM-based Methods |
|----------------------------|---|--|
| General description | Count the number of articles containing predefined risk-related keywords; measure the frequency of risk-related terms in text | Use attention-based neural networks to capture context, meaning, and relationships at the sentence and document level |
| Examples | EPU (Baker et al., 2016), FSS (Correa et al., 2021), GPR (Caldara & Iacoviello, 2022) | BERT (Devlin et al., 2019), FinBERT (Yang et al., 2020), CentralBankRoBERTa (Pfeifer et al., 2023), ChatGPT-based indicators (Bond et al., 2023 ; Kwon et al., 2025) |
| Advantages | Simple, transparent, interpretable, easy to implement | Can better capture context and forward-looking narratives, flexible, scalable to large volumes of text |
| Limitations | Lacks context (e.g. “tax cuts” may be positive or negative); performs poorly outside its original context | Higher complexity, limited interpretability, risk of hallucinations, more resource-intensive |

Recent work increasingly uses LLMs to complement traditional text-based indicators because they capture richer and more forward-looking narratives.

2

Method and Data

SPOT: Severity and Probability of Trigger events

- Goal: Develop a forward-looking indicator measuring the probability and severity of potential trigger events in the euro area (impact on euro area economic activity and/or financial stability)
- Method: Text-based identification and classification of potential trigger events using LLMs and Financial Times newspaper articles
- Classification: Articles are flagged if they signal a potential trigger event with severe macro-financial implications within the next year.
- Extracted attributes:
 - Probability that the trigger event materialises (low → high)
 - Severity (mild → very severe)
 - Time horizon (1–3m, 3–6m, 6–12m, longer)
 - Trigger source (macroeconomic, financial, geopolitical, monetary policy, fiscal, other)
- Indicator construction: trigger frequency, average probability, average severity and expected impact ($probability \times severity$), aggregated across all articles classified as triggers
- Why it matters: the four attributes allow decompositions, conditioning, and interpretation of risks — making the indicator more operational for analysis and monitoring compared to dictionary-based methods

Operationalising SPOT: Prompting pipeline

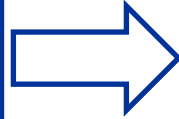
Step 1

Input: 1.1 million articles

Task: Pre-filter articles with an economic or financial stability focus

Output: 650k articles

Model: GPT-4o-mini



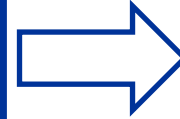
Step 2

Input: 650k articles

Task: Identify articles signalling severe macro-financial trigger events within the next year

Output: 140k articles

Model: GPT-4o-mini



Step 3

Input: 140k articles

Task: Extract risk attributes (e.g. severity, source, probability) for all articles flagged in Step 2

Output: 140k articles

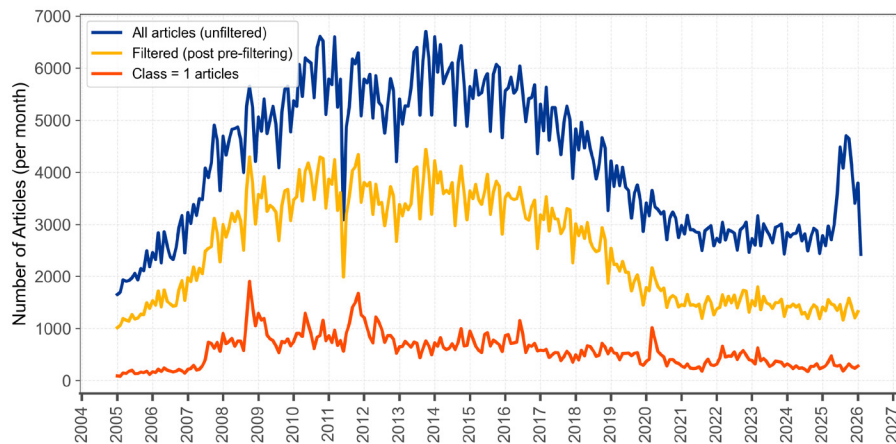
Model: GPT-4o-mini

In the final sample, 650k articles are retained, of which 140k articles are classified as signalling potential trigger events and assigned risk attributes.

A 20-year news dataset for financial stability analysis

- Full dataset covers a 20-year timespan of English articles published in the Financial Times (FT)
- More than 1 million articles are included, but unevenly distributed across months and years; pre-filtering removes some noise
- Articles published in the FT are quite heterogeneous

Number of FT articles varies strongly over time
(# articles per month)



FT article coverage varies substantially over time and content type (Descriptive statistics)

| Criterion | Mean | Std. Deviation | Min | Max |
|--------------------|-------|----------------|-------|-------|
| Words per article | 570 | 473 | 154 | 40476 |
| Articles per month | 4113 | 1463 | 1610 | 6705 |
| Articles per year | 44558 | 21801 | 22460 | 70663 |

Most articles were published on a Friday (20.7%), fewest on a Saturday (3.0%)

Results remain robust across sensitivity checks

- Stable results across different prompt specifications
- Limited sensitivity to model size and capacity
- Robustness to variations in sample size
- We explicitly mitigate potential information leakage
 - The LLM is instructed to avoid using external or future knowledge beyond the scope of the article
 - *(Prompt excerpt): "Use only information contained in the article as of its publication date, do not use outside knowledge, future information, or speculation."*
 - We do not ask the LLM to generate the aggregate indicator directly, but instead use it to label individual articles
 - Indicator performance does not appear to change when using articles published after the model training cut-off date

Example article — 29 January 2026

“Europe and the UK would lose more than the US in a trade war, research finds” (headline)

| Model | Probability | Severity | Horizon | Source |
|---------|-------------|----------|---------|---------------|
| 4o | 2 | 2 | 3 | Geo-political |
| 4o-mini | 2 | 2 | 3 | Geo-political |
| 5-mini | 2 | 2 | 2 | Geo-political |

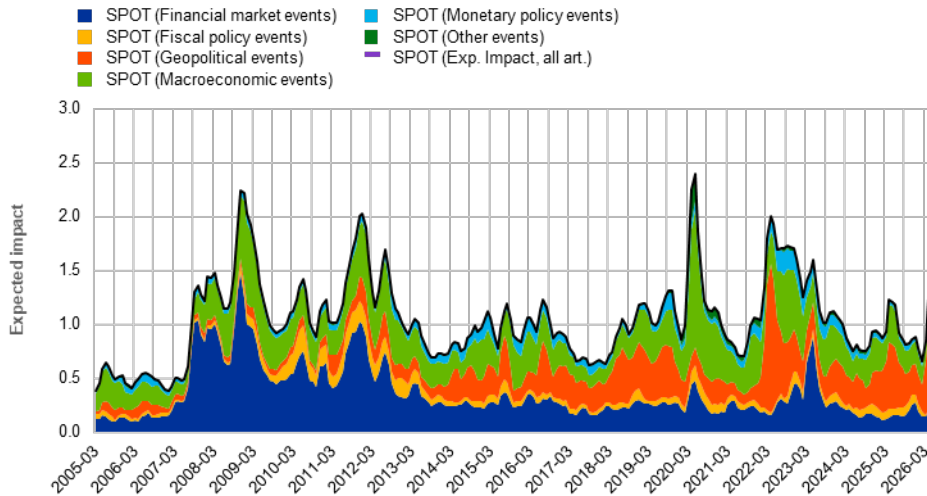
3

Results

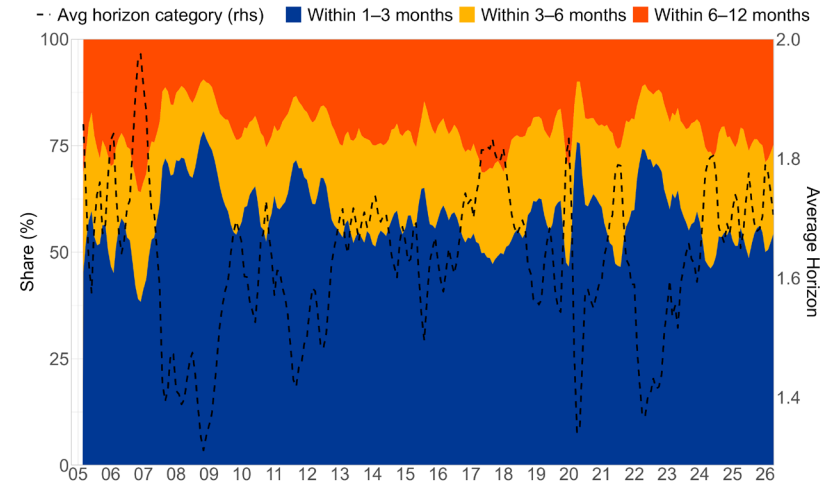
Different SPOT dimensions provide complementary monitoring signals

- SPOT indicators align with major historical stress episodes, including the GFC, EA debt crisis, COVID-19 pandemic and the Ukraine war
- Varying trigger sources: financial shocks in 2008–12, exogenous shock in 2020, geopolitical risk in 2022
- The share of short-term trigger articles increases ahead of crises, while trigger probabilities remain elevated at the current end

Decompositions into trigger sources help form a narrative (Decomposition of expected impact (probability x severity) across all articles)



Horizon of triggers varies over time (Share of trigger articles with different horizons of impact)

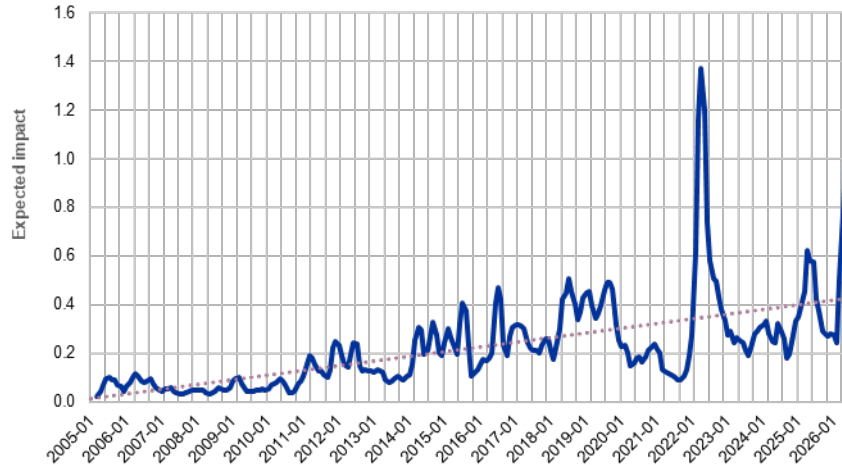


Sources: Financial Times data, authors' calculations.
Notes: Expected impact = average probability × severity across all articles, where non-trigger articles are assigned zero values.

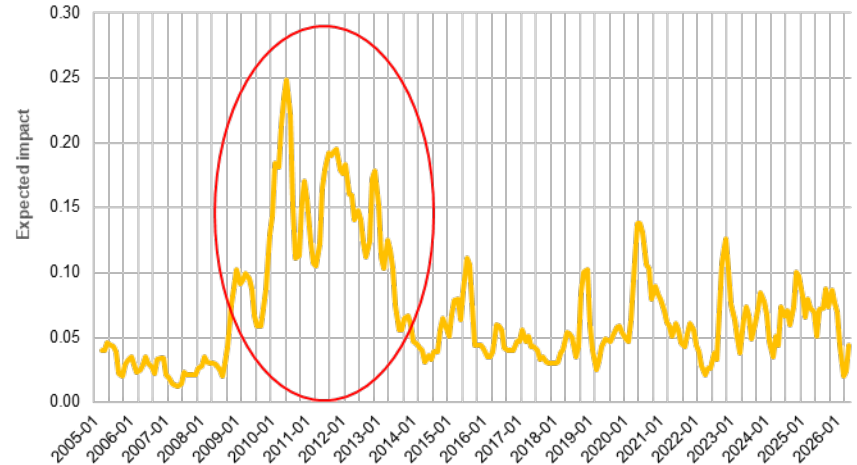
Sectoral SPOT indices can be used to monitor specific trigger sources

- Sectoral indices align well with past risk narratives
- Geopolitical triggers spiked in 2022 and 2025 and remain elevated relative to historical averages
- Fiscal policy triggers peaked during the euro area debt crisis and have risen again during 2025

Geopolitical triggers on an upward trend in past 15 years (Expected impact (Probability*Severity) for geopolitical trigger articles)



Fiscal policy triggers important during EA debt crisis (Expected impact (Probability*Severity) for fiscal policy trigger articles)

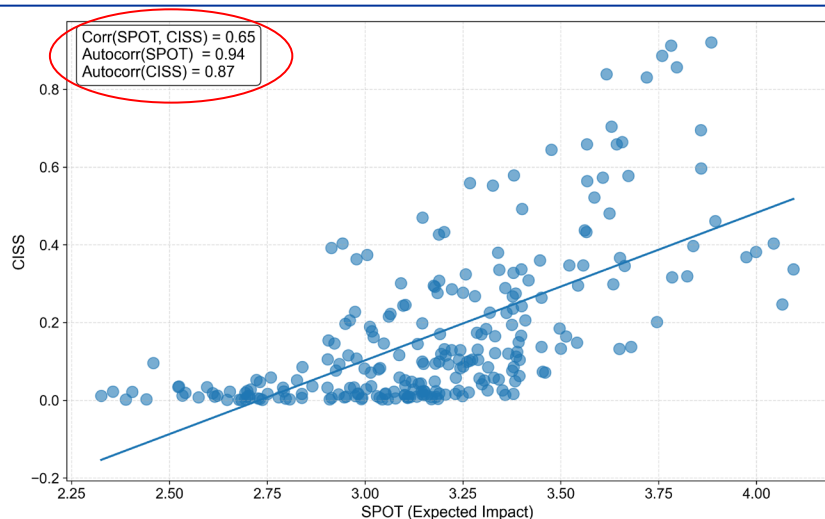


For comparison we also benchmark against other stress/risk indicators

- SPOT correlates with CISS (~ 0.7) but also contains complementary news-based information, for example, it is more persistent than the CISS
- SPOT captures additional dimensions relative to standard stress indicators and also accounts for the potential impact on the economy or financial system

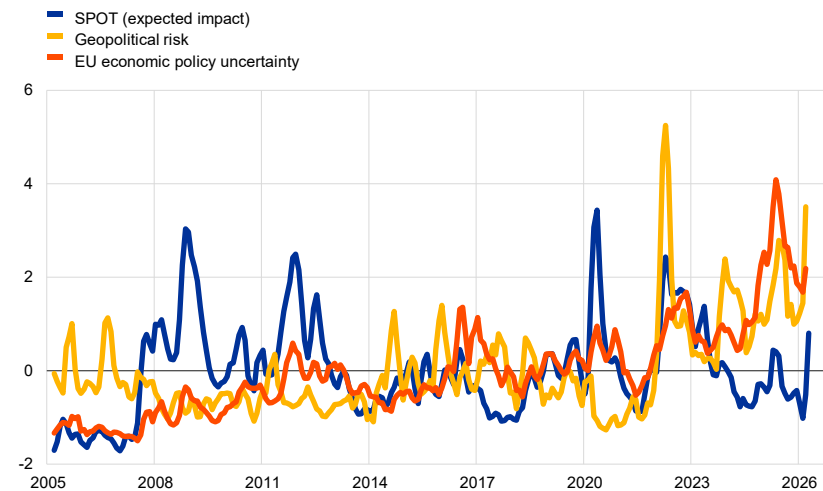
SPOT complements stress indicators like the CISS

(CISS: normalized, SPOT: Probability x Severity, Class = 1)



SPOT captures more dimensions than other indicators

(SPOT, geopolitical risk, and policy uncertainty; standardized indices)



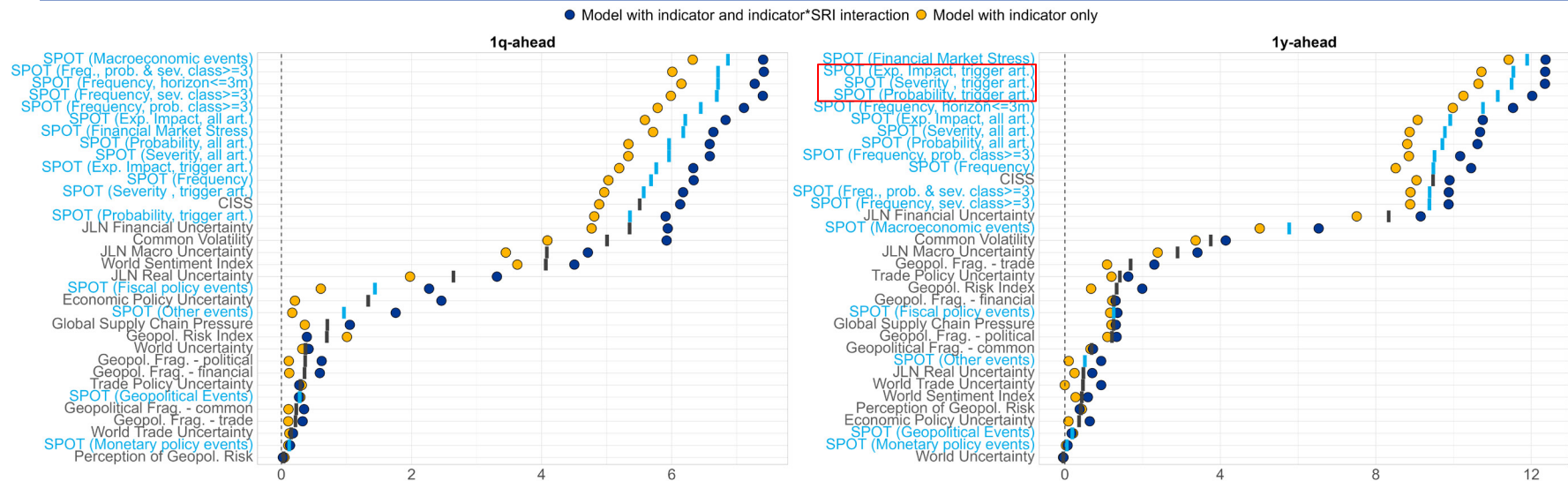
Sources: Financial times data, authors' calculations.

SPOT indicators outperform CISS and other indicators in G@R models

- A Growth-at-Risk (G@R) panel model is used to assess the information content of SPOT indicators
- SPOT indicators outperform other indicators for short-term (1q) and medium-term (1y) horizons
- The interaction between triggers and vulnerabilities appears important in driving downside risks to growth

Trigger event information improves the early identification of downside risks to the economy

(% improvement in panel growth-at-risk model fit (tick loss) for the 10th percentile at 1q- and 1y-ahead horizons)



Sources: Authors' calculations.

Notes: The baseline panel model includes GDP and the Systemic Risk Indicator (including its lag and a dummy variable for cases when SRI>0) and covid period dummies (2020q1-2020q4,2021q2). We use 3-month moving averages to convert higher frequency indicators to quarterly values.

4

Conclusion

Conclusion and outlook

- We develop **SPOT**, a forward-looking, news-based indicator capturing the probability and severity of financial stability trigger events
- **SPOT aligns well with major historical crisis episodes**, and the rich set of extracted attributes (probability, severity, horizon, source) helps form a risk narrative
- **SPOT contains complementary information** to existing risk measures (e.g. CISS, GPR) and improves short- and medium-term **Growth-at-Risk models**
- **SPOT can be integrated** into regular monitoring frameworks to complement vulnerability and stress indicators
- **Ongoing robustness assessments** will remain important as underlying models evolve (e.g. varying prompts, LLM vintages, and new data/information)
- **Working Paper and Financial Stability Review (May 2026) article** available

A SPOT in the dark: Using AI to assess financial stability risks

Domenic Kellner

European Central Bank, DG/MF
domenic.kellner@ecb.europa.eu

Jan Hannes Lang

European Central Bank, DG/MF
jan-hannes.lang@ecb.europa.eu

Lukas Joseph Nagy

European Central Bank, DG/MF
lukas_joseph.nagy@ecb.europa.eu

Marek Rusnák

European Central Bank, DG/MF
marek.rusnak@ecb.europa.eu

Acknowledgements

We would like to thank participants at National Bank of Slovakia Financial Stability Seminar, the IMF MCM Policy Forum, and the ECB Financial Stability Seminar for helpful comments and suggestions.

Disclaimer

The views expressed in this paper are those of the authors and do not necessarily reflect those of the European Central Bank.

Abstract

Financial stability risks are made up of two distinct components: vulnerabilities and possible trigger events. While there has been considerable progress regarding the measurement of vulnerabilities, the assessment of possible trigger events remains largely qualitative. To fill this gap, we employ Large Language Models to extract information about the **Severity and Probability Of potential Trigger events (SPOT)** from a large dataset of financial news articles over the period 2005 – 2025. Our SPOT indicator increases ahead of major historical trigger events, correctly identifies trigger sources, and helps to improve forward looking model estimates of downside risks to the economy. Our results indicate that the use of AI-based signal extraction from text can be a promising avenue to improve the monitoring of financial stability risks.

JEL codes

JEL: C55, C88, E32, E44, G01

Keywords

Financial stability; Artificial intelligence; Crisis indicators; Growth-at-risk;

Contents

- Non-technical summary** **1**

- 1 Introduction** **2**

- 2 Literature review** **4**

- 3 SPOT methodology and data** **6**
 - 3.1 Prompting strategy 6
 - 3.2 Dataset 8
 - 3.3 Construction of SPOT indicators 12

- 4 Empirical results** **13**
 - 4.1 Aggregate SPOT indicators and decompositions 14
 - 4.2 Granular SPOT indicators 18

- 5 Evaluation of SPOT indicators** **22**
 - 5.1 Comparison to other risk indicators from the literature 22
 - 5.2 SPOT performance evaluation in growth-at-risk models 28
 - 5.3 Robustness and sensitivity analysis of SPOT 31

- 6 Conclusion** **33**

- References** **34**

- Annex** **38**
 - A.1 Additional data and indicator characteristics 38
 - A.2 Robustness of Growth-at-Risk results 38

Non-technical summary

Financial stability risks typically arise from the interaction between underlying macro-financial vulnerabilities and adverse trigger events that can materialise unexpectedly. While substantial progress has been made in monitoring vulnerabilities, the systematic assessment of potential trigger events remains challenging and often relies on qualitative judgement. This paper proposes a novel framework to quantify such trigger events using artificial intelligence (AI) and large-scale text analysis.

We introduce SPOT, a novel AI-based indicator to measure the **Severity and Probability Of Triggers**. SPOT is generated using large language models (LLMs) and a large dataset of financial news articles. The LLM assesses whether an article describes events capable of adversely affecting euro area economic activity or financial stability in the future and then performs a structured assessment along four dimensions: the probability that the event could materialise, the severity of its potential macro-financial impact, the expected time horizon of the impact, and the dominant source of the trigger. These article-level assessments are then aggregated into the SPOT indicator to capture both the frequency and intensity of potential triggers over time.

The resulting SPOT indicator tracks major episodes of financial and macroeconomic stress and provides detailed information on the nature of emerging risks. The benchmark SPOT indicator rises ahead of the global financial crisis, the euro area sovereign debt crisis, the COVID-19 pandemic and periods of heightened geopolitical tensions. The framework also allows for a decomposition of SPOT by trigger source and by attributes such as probability, severity and expected timing, thereby offering a structured and comprehensive view of evolving financial stability risk perceptions and their underlying drivers.

We further show that SPOT provides information beyond existing risk indicators, including measures of financial stress, geopolitical risk and economic policy uncertainty. While correlated with these indicators, SPOT captures broader forward-looking narratives reflected in financial news and incorporates information about the probability and potential impact of adverse events. This additional information proves economically meaningful: incorporating SPOT into a growth-at-risk framework significantly improves the assessment of downside risks to economic growth, particularly when combined with measures of financial vulnerabilities.

Overall, the results demonstrate that AI-based text analysis can enhance financial stability monitoring by systematically extracting forward-looking signals from large volumes of unstructured information.

1 Introduction

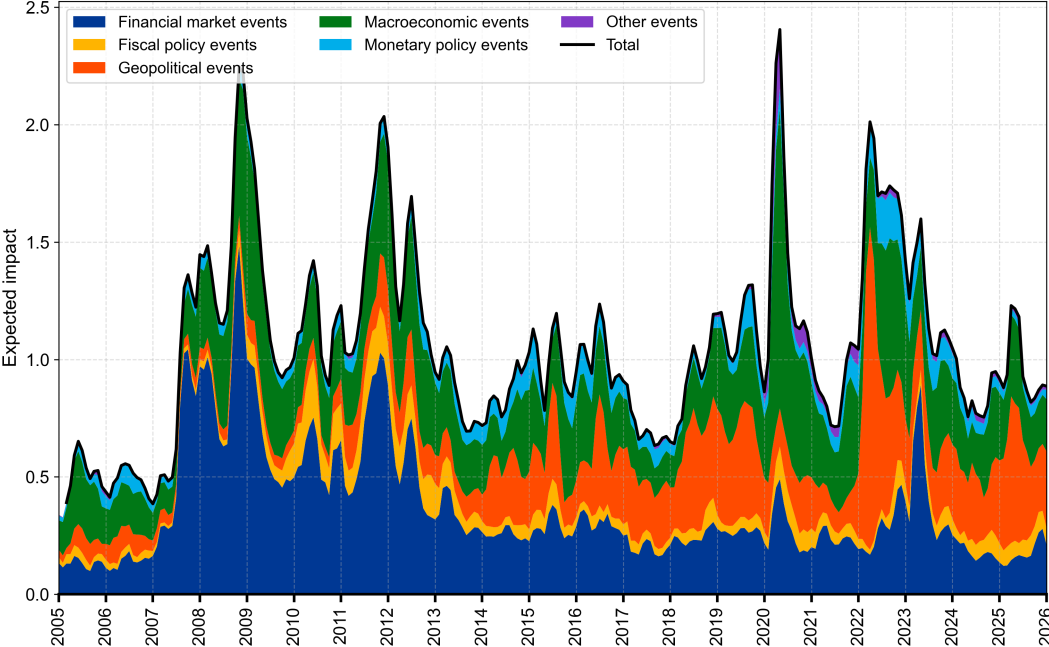
Since the global financial crisis, the analysis of financial stability risks has gained prominence at central banks, supervisory authorities, and international institutions. At a conceptual level, financial stability risks can be broken down into two components: 1) vulnerabilities, such as excessive leverage, maturity mismatches, or asset price misalignments; and 2) potential trigger events, such as bank failures, interest rate shocks, a loss of confidence, wars, or pandemics (Covitz et al., 2015; Fell and Schinasi, 2005; Fell et al., 2024). While there has been considerable progress in measuring vulnerabilities using individual or composite indicators (Alessi and Detken, 2011; Borio and Drehmann, 2009; Borio and Lowe, 2002; Detken et al., 2014; Lang et al., 2019; Schularick and Taylor, 2012; Schüler et al., 2020), the monitoring of potential triggers remains largely qualitative due to the complex and constantly evolving nature of such triggers. Hence, the assessment of financial stability risks and especially of potential triggers is akin to walking barefoot in a dark room filled with broken glass: you are not sure whether the next step is risky or not, due to limited visibility regarding what is lying ahead.

The aim of this paper is to generate an additional SPOT of light into the analysis of financial stability risks, by employing recent advances in Artificial Intelligence (AI) to systematically assess the potential Severity and Probability Of Triggers (SPOT) based on news. Specifically, we apply a structured Large Language Model (LLM) prompt to more than 1 million newspaper articles from the *Financial Times* starting in 2005 to classify whether an article signals a potential future trigger event that could have a severe negative impact on the euro area economy and financial stability, and to assess various properties of the potential trigger, such as the severity, probability, time horizon, and trigger source. The motivation for this approach is that newspaper articles often discuss potential trigger events and LLMs can be used to extract this information systematically. Based on the article classification and extracted trigger properties from our structured prompting process, we construct various SPOT indicators that can be used to complement vulnerability indicators for the assessment of financial stability risks.

Our benchmark SPOT indicator, which reflects the average expected impact of trigger events across all articles, spikes ahead of major historical episodes of financial instability, such as the global financial crisis, the euro area sovereign debt crisis, the Covid-19 pandemic, or the Russian invasion of Ukraine (See Figure 1). Moreover, the decomposition of the SPOT indicator by trigger source allows for further insights into underlying drivers and helps form a risk narrative. For instance, financial market triggers dominated during the global financial crisis, other exogenous triggers increased in importance during the COVID-19 pandemic, while geopolitical triggers spiked following Russia's invasion of Ukraine in 2022. Moreover, sectoral SPOT

indicators can be used to complement the benchmark SPOT indicator for monitoring of specific trigger sources, such as those related to geopolitical events.

Figure 1: Benchmark SPOT indicator with decomposition by trigger source over time



Notes: The benchmark SPOT indicator is calculated as the average expected impact (probability*severity) of trigger events across all articles at a given point in time.

We also use the growth-at-risk modeling approach pioneered by [Adrian et al. \(2019\)](#) to formally evaluate our AI-based indicators. The estimation results from euro area panel growth-at-risk models for 1-quarter and 1-year ahead horizons show that various SPOT indicators contain useful information about future downside risks to the economy and outperform other commonly used indicators capturing financial stress ([Kremer et al., 2012](#)), geopolitical risk ([Caldara and Iacoviello, 2022](#)), policy uncertainty ([Baker et al., 2016](#)) or volatility ([Engle and Campos-Martins, 2023](#)). The estimation results also reveal that interacting SPOT indicators with commonly used vulnerability indicators in growth-at-risk models ([Lang et al., 2025](#)) further improves the explanatory power of the models, indicating important amplification mechanisms between vulnerabilities and triggers in driving economic tail risks. Overall, the results presented in this paper suggest that the AI-based SPOT indicators can measure potential trigger events and underlying trigger sources in a systematic way and provide useful information for the monitoring of financial stability risks.

Our paper contributes to the literature on the use of AI to generate indicators from text, the

measurement of financial stability risks, and the use of growth-at-risk models. The remainder of the paper is structured as follows. Section 2 discusses the related literature in more detail. Section 3 presents our prompting strategy, the data used, and the proposed SPOT indicators. Section 4 provides an overview of empirical results for various SPOT indicators. In Section 5 we compare the SPOT indicators to other indicators suggested in the literature, perform the panel growth-at-risk model estimations, and show various robustness exercises. Section 6 concludes.

2 Literature review

Textual data potentially contain economically relevant information that can complement traditional numeric data sources. Methods for extracting information from text typically require processing large corpora (Gentzkow et al., 2019) and have benefited substantially from recent advances in computational power and language modeling. Our paper relates to a growing literature that uses text-based methods to construct economic and financial indicators, and in particular to recent work employing Large Language Models (LLMs) for sentiment measurement, uncertainty quantification, and financial stability analysis.

A large body of earlier research constructs economic indicators from text using relatively simple Natural Language Processing (NLP) techniques such as keyword dictionaries, word counts, or bag-of-words approaches. Prominent examples include the Economic Policy Uncertainty (EPU) index by Baker et al. (2016), the Financial Stability Sentiment (FSS) index by Correa et al. (2017), and the Geopolitical Risk (GPR) index by Caldara and Iacoviello (2022). While these methods are transparent and easy to implement, they rely on predefined vocabularies and typically fail to capture semantic context, narrative structure, or evolving economic language (Loughran and McDonald, 2011). Recent work therefore seeks to improve text-based measurement by incorporating contextual language models and neural-network-based classifiers.

Several studies demonstrate the potential of LLMs for generating sentiment-based indicators from financial news. Bond et al. (2023) use ChatGPT to construct sentiment-based market indicators from daily news summaries. Lopez-Lira and Tang (2023) show that GPT-based models can extract economically meaningful signals from news headlines that predict market reactions. Similarly, Lefort et al. (2024) apply LLMs to financial news headlines to derive sentiment scores for NASDAQ index predictions. Zhang et al. (2025) propose *FinSentLLM*, which combines an ensemble of LLMs with structured financial signals and expert-driven semantic cues to improve financial sentiment forecasting. Their approach captures inter-model agreement and domain-

specific information and yields improvements in predictive accuracy and long-run comovement with major stock indices.

Beyond sentiment scores, [Nyman et al. \(2021\)](#) show that changes in the coherence and emotional structure of financial narratives systematically precede episodes of financial instability. Their findings support the view that news text embeds forward-looking information about systemic risk, providing a conceptual foundation for event-based monitoring frameworks. While this work provides an important conceptual foundation, monitoring of macro-financial risks often requires more compressed and easily interpretable indicators. [Kwon et al. \(2025\)](#) construct sentiment indices for macroeconomic growth and inflation from a large corpus of U.S. news articles. Their approach directly extracts economic narratives from press coverage and decomposes sentiment into interpretable drivers such as demand and supply forces. The resulting indicators closely track conventional macroeconomic measures and improve forecasting performance when incorporated into predictive models, suggesting that text-based sentiment captures information not contained in traditional data sources.

A related strand of literature applies LLMs to measure economic uncertainty and policy-related risks. [Audrino et al. \(2024\)](#) introduce an LLM-based framework for measuring economic uncertainty from newspaper texts that exploits contextual language understanding to classify sources of uncertainty across domains such as geopolitics, economic policy, monetary policy, and financial markets. The resulting indices exhibit stronger correlations with macroeconomic developments and financial market volatility than traditional benchmarks. Similarly, [Chen et al. \(2024\)](#) revisit the construction of the Economic Policy Uncertainty index by showing that keyword-matching approaches neglect semantic context and introduce noise. By replacing keyword matching with neural-network classifiers, they improve predictive performance and out-of-sample accuracy.

Our paper contributes to this literature by proposing a structured LLM-based framework to identify potential financial stability trigger events from news in a forward looking manner and to assess their probability, severity, timing, and source. While existing work primarily focuses on sentiment or uncertainty measurement using LLMs, we construct SPOT indicators that explicitly capture potential trigger events for financial stability risks and examine their interaction with financial vulnerabilities within a growth-at-risk framework.

3 SPOT methodology and data

Financial stability risks arise from a combination of underlying vulnerabilities and potential trigger events (Covitz et al., 2015; Fell and Schinasi, 2005; Fell et al., 2024), yet such triggers remain inherently difficult to quantify due to their complex and evolving nature. To address this gap, we develop a forward-looking SPOT indicator to systematically identify and characterise the Severity and Probability Of Triggers (SPOT) for financial stability risks in the euro area. The SPOT indicator is constructed from *Financial Times* (FT) articles evaluated using GPT-4o-mini within a structured multi-stage prompting pipeline designed for large-scale text classification.

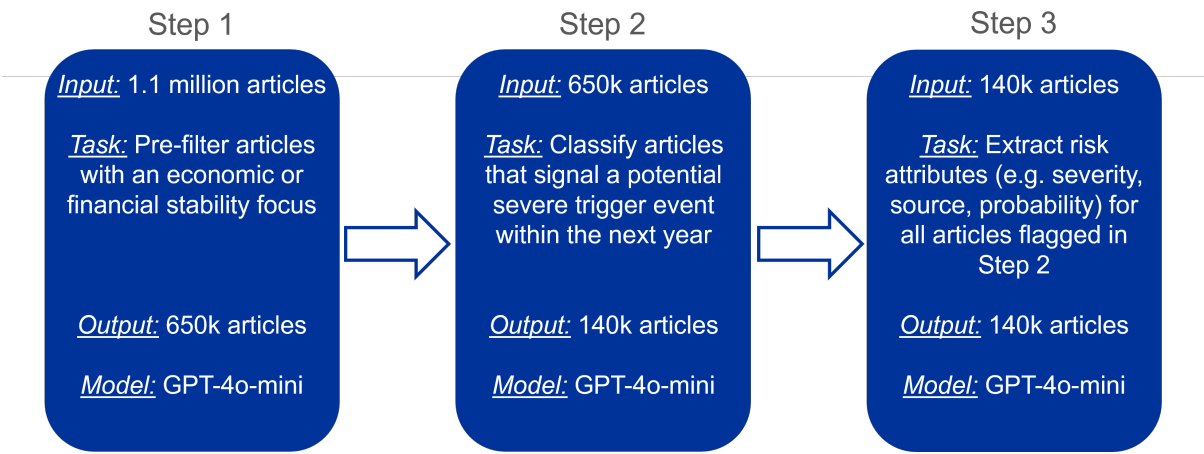
3.1 Prompting strategy

The proposed prompting framework decomposes the task into three sequential steps (Figure 2) to first identify economically relevant content, then detect forward-looking trigger events, and finally extract key attributes of potential trigger events. Each step represents a different prompt that we ask the LLM to evaluate via an API. This modular prompt design mitigates context window limitations, reduces cognitive load on the model, and improves the consistency and interpretability of the resulting classifications. The prompting framework relies on rigorously specified instructions that define key economic and financial stability concepts, enforce forward-looking discipline, and structure model outputs. In addition, the model is explicitly instructed to use only information contained in the article as of its publication date and to avoid using external knowledge, future information, or speculation. This restriction mitigates information contamination and helps to ensure that the resulting classifications reflect real-time information contained in the newspaper article, consistent with the risk monitoring objective.

In the first step, all news articles are pre-filtered to identify content whose primary focus is economic or financial and whose economic or financial discussion is explicit, substantive, and central to the article. Economic and financial content is defined to include macroeconomic variables (such as inflation, GDP, interest rates, or unemployment), financial markets or instruments (including equities, bonds, spreads, banks, or non-bank financial institutions), or economic policy with a direct and concrete economic impact. Articles concerning individual firms are retained only if firm-specific developments are described as having implications for the broader economy or financial system. Articles that are strictly non-economic or non-financial, or that use economic terminology only metaphorically or without analytical substance, are excluded. When the classification is ambiguous, the model is instructed to conservatively exclude the article. This pre-filtering step removes unrelated content (e.g. articles only covering politics, art,

travel, etc.) and ensures that subsequent analysis focuses on articles with relevant information.

Figure 2: Three-stage prompting pipeline



Notes: Each stage represents a different prompt that is sent to the LLM to evaluate via an API.

In the second step, the pre-filtered articles are classified whether they signal current trigger events or potential future trigger events that could have a severe negative impact on economic activity or the stability of the financial system within the euro area over a 1–12 month horizon (*Class* = 1). A severe negative impact on economic activity is defined as negative real GDP growth, significant declines in industrial production, or substantial increases in unemployment, while a severe negative impact on financial stability includes banking-sector losses, significant increases in interbank lending rates, or large outflows of bank deposits. The classification explicitly adopts a forward-looking perspective, identifying statements about forecasts, expectations, warnings, or predictions concerning potential future developments affecting the euro area. Articles describing only current or past developments without credible future implications for euro area economic activity or financial stability are classified as non-trigger events. Particular emphasis is placed on euro area relevance: articles concerning non-euro area developments or firm-specific events are classified as trigger events only if plausible spillovers to euro area economic activity or financial stability within the specified horizon are indicated or strongly implied.

In the third step, articles identified as signaling potential trigger events are further characterised along several dimensions capturing the probability (*Prob*), severity (*Sev*), time horizon, and source of the potential trigger. To provide an additional consistency check for the second stage, the third stage explicitly allows zero-value categories, thereby giving the model an exit option even when an article has been previously classified as signaling a trigger event. This

design acts as a form of double verification, allowing the model to revise earlier assessments if the evidence for a genuine macro-financial risk trigger is weak upon closer evaluation.

Probability is measured on a four-point scale ranging from 0 (unlikely, implying effectively zero probability) to 3 (high probability, indicating a very plausible outcome likely to materialise), with intermediate values representing low and medium probability. Severity is also measured on a four-point scale, where 0 indicates negligible or limited effects on output, employment, or financial stability; 1 indicates mildly severe outcomes such as contained sectoral stress or limited financial tightening; 2 indicates severe outcomes involving broad-based macroeconomic or financial disruption; and 3 indicates very severe or systemic impacts, including banking distress, sovereign funding crises, or widespread economic contraction.

In the third step, the model also assigns a time horizon indicating when the impact is most likely to materialise, distinguishing between impacts that have already materialised, impacts expected within 1–3 months, 3–6 months, 6–12 months, and 12–36 months. Finally, the model identifies the dominant source of the trigger event by assigning the article to one of several predefined categories. These include macroeconomic events (such as demand disturbances, commodity price shocks, or exchange rate adjustments), financial market events (including liquidity shortages, banking-sector instability, or reassessments of asset risk), geopolitical events (such as wars, political instability, trade conflicts, or sanctions), monetary policy events, fiscal policy events, or other (exogenous) shocks such as natural disasters or pandemics.

Overall, the three-step prompting pipeline enables scalable and systematic identification of potential trigger events from a large set of newspaper articles, while maintaining conceptual clarity. The modular structure further enhances operational flexibility, as individual steps can be updated independently in case one wants definitions or emerging trigger source categories to be adjusted. In particular, revisions to the attribute characterizations can be implemented by modifying only the final prompting stage without re-running the entire pipeline, which facilitates continuous monitoring and regular updating of risk assessments. The sequential prompting design also improves computational efficiency by restricting more demanding classification tasks to a progressively smaller subset of relevant articles.

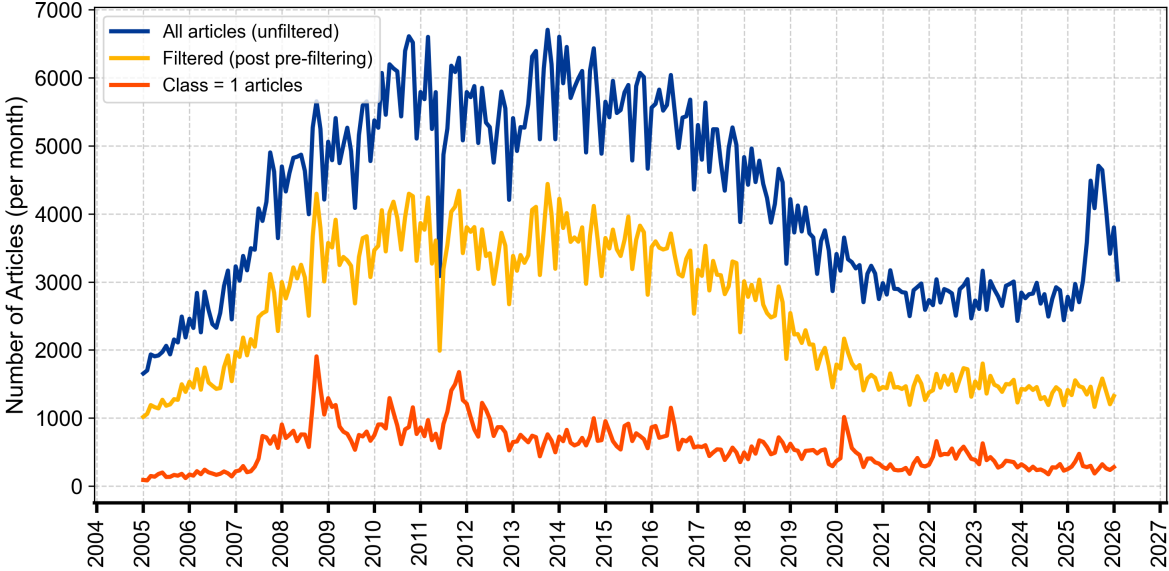
3.2 Dataset

We apply the three-stage prompting framework described above to a large dataset of newspaper articles from the *Financial Times* (FT). The FT serves as our data source due to its consistent high-quality focus on economic and financial topics in Europe, which ensures thematic

coherence and reduces the need for extensive pre-screening of unrelated content. The dataset covers articles published between January 2005 and January 2026 and provides broad coverage of macroeconomic and financial developments relevant for financial stability analysis. For each article, we extract the title, body text, and publication date, which serve as inputs to the LLM classification procedure. The publication date ensures that model assessments rely only on information available at the time of publication (as instructed in the LLM system prompt), thereby supporting the forward-looking design of the framework and mitigating potential look-ahead bias.

The raw dataset comprises approximately one million articles, with an average length of around 600 words and roughly 4,000 articles published per month (Figure 3). Overall data quality is high. Articles with missing body text or incomplete metadata are excluded, and duplicate entries are removed from the sample. The number of available articles varies over time, reflecting changes in archiving practices, topic coverage, and distribution rights that affect data availability. This variation underscores the importance of constructing indicators based on article content rather than relying solely on article counts, a potential limitation of many conventional dictionary-based approaches.

Figure 3: Number of Articles over time for the three prompting stages



Notes: The blue line shows the number of articles per month before the first stage i.e., the raw data. The yellow line shows the number of articles before the second stage and after pre-filtering articles for their economic relevance. The orange line shows the number of articles per month classified as *Class1* in the second stage, that are passed to the stage three prompt.

The multi-stage prompting procedure sequentially reduces the number of articles that need

to be processed in each step (Figure 2). After initial data cleaning, approximately 1.05 million articles enter the first stage of the pipeline, where the pre-filter removes articles without a primary economic or financial focus (around 39% of the sample). The second stage evaluates the remaining articles for forward-looking trigger events relevant for euro area economic activity or financial stability and excludes a further 78% of observations. In total, approximately 144,000 articles are retained for the final stage, in which potential trigger events are characterised along several dimensions such as probability, severity, time horizon and trigger source (Figure 3).

Table 1: Marginal distributions of extracted trigger attributes (Class 1 articles)

| Attribute | Category | Total (N) | Share (%) |
|----------------|---|-----------|-----------|
| Trigger source | Financial market events | 50,989 | 35.4 |
| | Macroeconomic events | 44,717 | 31.0 |
| | Geopolitical events | 25,924 | 18.0 |
| | Fiscal policy events | 10,915 | 7.6 |
| | Monetary policy events | 9,583 | 6.6 |
| | Other events | 1,981 | 1.4 |
| Probability | Unlikely / effectively zero probability (0) | 43 | 0.0 |
| | Low probability (1) | 4,221 | 2.9 |
| | Medium probability (2) | 93,881 | 65.1 |
| | High probability (3) | 45,964 | 31.9 |
| Severity | Negligible / limited effects (0) | 74 | 0.1 |
| | Mildly severe (contained stress) (1) | 10,678 | 7.4 |
| | Severe (broad disruption) (2) | 113,309 | 78.6 |
| | Very severe / systemic (3) | 20,048 | 13.9 |
| Horizon | Already materialising (0) | 5,428 | 3.8 |
| | Within 1–3 months (1) | 84,060 | 58.3 |
| | Within 3–6 months (2) | 27,107 | 18.8 |
| | Within 6–12 months (3) | 27,241 | 18.9 |
| | Within 12–36 months (4) | 273 | 0.2 |

Notes: Shares are computed relative to all Class 1 articles. Probability and severity are measured on a 0–3 scale and horizon on a 0–4 scale.

The marginal distributions of the extracted attributes from the third stage prompt are reported in Table 1. Most trigger articles are assigned a medium probability and severity, while high-probability and high-severity events occur less frequently, reflecting the rarity of crisis level triggers in news coverage. As a robustness feature of the prompting design, in the third stage the model is allowed to assign category 0 for probability and severity even for articles previously identified as trigger events. Such cases are rare ($\leq 0.1\%$), confirming the validity of the step 2 classification. For the horizon attribute, category 0 indicates that the adverse impact is already materialising rather than expected in the future. This occurs in only 3.8% of trigger articles and is most likely due to the fact that during crisis periods the distinction between triggers being

contemporaneous rather than forward-looking becomes somewhat blurry. More than half of the potential triggers are assessed to have an impact within 1-3 months. In terms of trigger sources, around 30% of triggers relate to financial market and macroeconomic events respectively, with a further 20% related to geopolitical events.

The joint distributions of extracted trigger attributes in Table 2 provide further insights. Slightly less than half of the high probability triggers are also high severity triggers. On the other hand, almost all of the high severity triggers also have a high probability of materialising. Around 90% of the high probability and high severity triggers are associated with short horizons. The occurrence of high probability and high severity triggers is therefore concentrated during periods of imminent stress. Conversely, the model rarely assigns high probability or high severity to triggers with distant horizons (≥ 3), which is consistent with the notion that news coverage of distant triggers will often be more speculative and vague when it comes to their potential impact. The distributions of severity, probability, and horizon characteristics are rather similar across the different trigger sources. However, short horizon triggers are dominated by financial market events, whereas medium to longer horizon triggers are dominated by macroeconomic events (Table 2).

Table 2: Joint distributions of trigger attributes (shares in % of Class 1 articles)

| Attribute | Category | Probability | | | | Severity | | | | Horizon | | | | Trigger source | | | | | | |
|----------------|----------|-------------|----|----|----|----------|----|----|----|---------|----|----|----|----------------|-------|--------|-----|----|--------|-------|
| | | 0 | 1 | 2 | 3 | 0 | 1 | 2 | 3 | 0 | 1 | 2 | 3 | 4 | Macro | FinMkt | Geo | MP | Fiscal | Other |
| Probability | 0 | | | | | 100 | 0 | 0 | 0 | 98 | 0 | 0 | 0 | 2 | 74 | 9 | 0 | 0 | 0 | 16 |
| | 1 | | | | | 1 | 98 | 1 | 0 | 1 | 22 | 17 | 57 | 4 | 42 | 27 | 11 | 6 | 11 | 3 |
| | 2 | | | | | 0 | 7 | 93 | 0 | 0 | 49 | 25 | 25 | 0 | 34 | 31 | 18 | 7 | 8 | 1 |
| | 3 | | | | | 0 | 0 | 57 | 43 | 11 | 80 | 7 | 2 | 0 | 23 | 45 | 19 | 6 | 6 | 1 |
| Severity | 0 | 58 | 42 | 0 | 0 | | | | | 59 | 0 | 0 | 14 | 27 | 73 | 9 | 0 | 5 | 0 | 12 |
| | 1 | 0 | 39 | 61 | 0 | | | | | 0 | 35 | 20 | 43 | 2 | 43 | 26 | 12 | 8 | 10 | 2 |
| | 2 | 0 | 0 | 77 | 23 | | | | | 1 | 58 | 21 | 20 | 0 | 32 | 33 | 18 | 7 | 8 | 1 |
| | 3 | 0 | 0 | 2 | 98 | | | | | 20 | 74 | 4 | 2 | 0 | 16 | 53 | 20 | 2 | 7 | 1 |
| Horizon | 0 | 1 | 1 | 3 | 96 | 1 | 1 | 24 | 75 | | | | | | 23 | 61 | 6 | 1 | 5 | 2 |
| | 1 | 0 | 1 | 55 | 44 | 0 | 4 | 78 | 18 | | | | | | 26 | 40 | 18 | 7 | 8 | 1 |
| | 2 | 0 | 3 | 86 | 11 | 0 | 8 | 89 | 3 | | | | | | 41 | 27 | 18 | 7 | 7 | 1 |
| | 3 | 0 | 9 | 88 | 4 | 0 | 17 | 82 | 1 | | | | | | 39 | 23 | 22 | 6 | 8 | 2 |
| | 4 | 0 | 68 | 22 | 10 | 7 | 65 | 22 | 5 | | | | | | 45 | 24 | 10 | 4 | 10 | 7 |
| Trigger source | Macro | 0 | 4 | 72 | 24 | 0 | 10 | 82 | 7 | 3 | 49 | 25 | 24 | 0 | | | | | | |
| | FinMkt | 0 | 2 | 57 | 41 | 0 | 5 | 74 | 21 | 7 | 67 | 15 | 12 | 0 | | | | | | |
| | Geo | 0 | 2 | 65 | 33 | 0 | 5 | 79 | 16 | 1 | 57 | 18 | 23 | 0 | | | | | | |
| | MP | 0 | 3 | 70 | 28 | 0 | 9 | 87 | 5 | 1 | 64 | 19 | 17 | 0 | | | | | | |
| | Fiscal | 0 | 4 | 70 | 26 | 0 | 9 | 78 | 12 | 3 | 59 | 17 | 21 | 0 | | | | | | |
| | Other | 0 | 7 | 67 | 26 | 0 | 11 | 75 | 14 | 7 | 54 | 14 | 24 | 1 | | | | | | |

Notes: Entries are row-wise shares in %. Within each row, reported blocks sum to 100. Probability is measured on a four-point scale from 0 (unlikely, effectively zero probability) to 3 (high probability). Severity is measured on a four-point scale from 0 (negligible impact) to 3 (very severe or systemic impact). The horizon indicates the expected timing of the impact: 0 = already materialising, 1 = within 1–3 months, 2 = within 3–6 months, 3 = within 6–12 months, and 4 = within 12–36 months. Trigger source categories are abbreviated as follows: Macro (macroeconomic events), FinMkt (financial market events), Geo (geopolitical events), MP (monetary policy events), Fiscal (fiscal policy events), and Other (exogenous events).

3.3 Construction of SPOT indicators

The classification into trigger articles ($Class = 1$) and non-trigger articles ($Class = 0$) from prompting step 2 and the extracted attributes from prompting step 3 can be used to create various SPOT indicators reflecting the potential **S**everity and **P**robability **O**f **T**rigger **S**ignals (**SPOT**). Before presenting these indicators, it is useful to define some notation. Let N_t denote the total number of economic or financial articles at a given point in time t , and let N_t^{C1} and N_t^{C0} denote the number of $Class1$ and $Class0$ articles respectively, where $N_t = N_t^{C1} + N_t^{C0}$. Whenever we sum over the entire set of articles we use the index i and whenever we sum over a subset of articles we use the index j . By default, whenever an article is classified as $Class0$ in prompting step 2, all of the attributes that are extracted in prompting step 3 are automatically set to zero. Based on these notation conventions we can easily define the following SPOT indicators.

Trigger Frequency Indicator. The share of articles signaling potential trigger events:

$$\overline{Freq}_t = \frac{1}{N_t} \sum_{i=1}^{N_t} Class_i = \frac{1}{N_t} \sum_{j=1}^{N_t^{C1}} Class_j = \frac{N_t^{C1}}{N_t} \quad (1)$$

Average Probability Indicator. The average probability of triggers across all articles:

$$\overline{Prob}_t = \frac{1}{N_t} \sum_{i=1}^{N_t} Prob_i = \frac{1}{N_t} \sum_{j=1}^{N_t^{C1}} Prob_j \quad (2)$$

Average Severity Indicator. The average severity of triggers across all articles:

$$\overline{Sev}_t = \frac{1}{N_t} \sum_{i=1}^{N_t} Sev_i = \frac{1}{N_t} \sum_{j=1}^{N_t^{C1}} Sev_j \quad (3)$$

Average Expected Impact Indicator (Benchmark SPOT indicator). This measure captures both the frequency and intensity of trigger signals:

$$\overline{Impact}_t = \frac{1}{N_t} \sum_{i=1}^{N_t} Prob_i \times Sev_i = \frac{1}{N_t} \sum_{j=1}^{N_t^{C1}} Prob_j \times Sev_j \quad (4)$$

Decompositions of Indicators. The SPOT indicators above can also be systematically decomposed by a certain attribute A (e.g. by trigger source, time horizon, probability, or severity), where N^A denotes the number of different attribute types and $N_t^{C1|a}$ denotes the number of trigger articles associated with a specific attribute type a . Such decompositions allow for a more detailed analysis of the nature of triggers and can help form a coherent risk narrative:¹

$$\overline{\text{Measure}}_t = \sum_{a=1}^{N^A} \left(\frac{1}{N_t} \sum_{j=1}^{N_t^{C1|a}} \text{Measure}_j \right) \quad (5)$$

Sectoral SPOT Indicators. Frequency, average probability, severity, and expected impact can also be calculated for different trigger sources S , where $N_t^{C1|S}$ denotes the number of trigger articles associated with a specific trigger source:

$$\overline{\text{Measure}}_t^S = \frac{1}{N_t} \sum_{j=1}^{N_t^{C1|S}} \text{Measure}_j \quad (6)$$

Trigger Article Indicators. Average probability, severity, and expected impact can also be calculated conditioning on trigger articles only:

$$\overline{\text{Measure}}_t^{C1} = \frac{1}{N_t^{C1}} \sum_{j=1}^{N_t^{C1}} \text{Measure}_j \quad (7)$$

The default time aggregation of articles into SPOT indicators is done at a monthly frequency. To smooth short-term fluctuations and enhance interpretability, all indicators are transformed into three-month moving averages.

4 Empirical results

This section provides an overview of the empirical properties of the various SPOT indicators described in [subsection 3.3](#).

¹ The decomposition of expected impact by probability/severity is not possible. For SPOT indicators representing the probability or severity, the decompositions can be conducted along the complementary dimension.

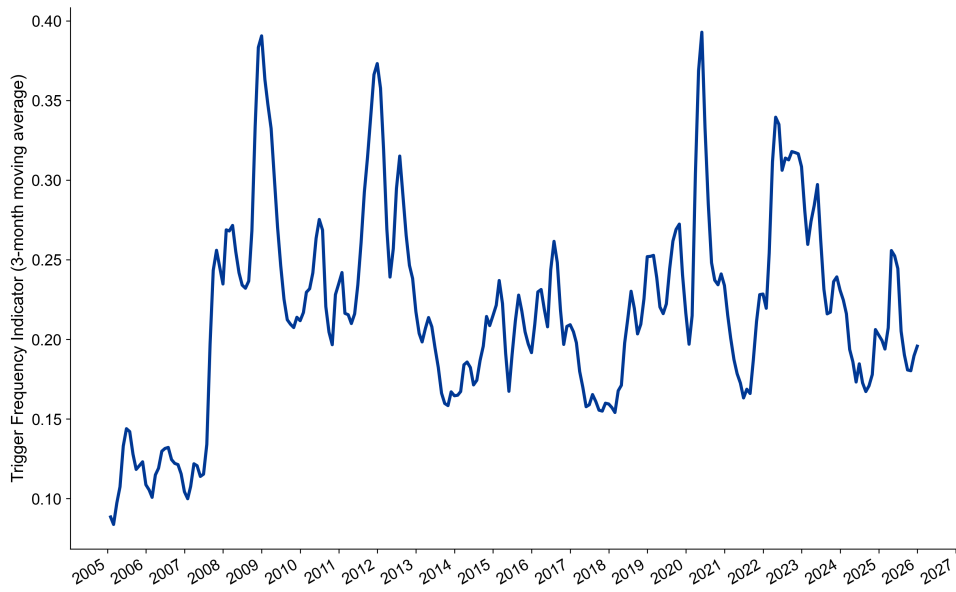
4.1 Aggregate SPOT indicators and decompositions

The first observation is that the frequency of trigger articles varies considerably over time, between 10% and close to 40% (Figure 4 panel a). The largest peaks in the frequency of trigger articles are reached during the global financial crisis in 2008, and the Covid-19 crisis in 2020. The second observation is that while the average probability and average severity of trigger articles display a high positive correlation, they also seem to contain complementary information (Figure 4 panel b). For example, the amplitude of the probability is almost twice as large as the amplitude of the severity. Moreover, there are time periods when these two attributes of potential triggers show different dynamics, for example during 2025. As shown in section 5.2 below, combining information about the frequency, probability, and severity of trigger articles improves model estimates of future downside risks to the economy. Hence, our benchmark SPOT indicator (the average expected impact) combines these three dimensions.

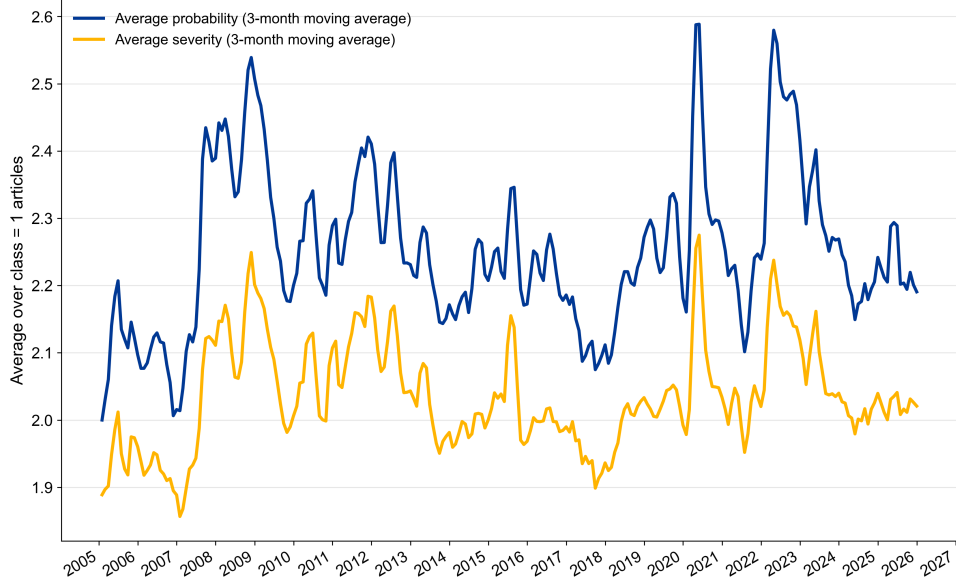
The benchmark AI-based SPOT indicator increases ahead of major historical trigger events and correctly identifies the relative importance of different trigger sources, as shown in Figure 5. For example, the benchmark SPOT indicator reaches peaks around the onset of the global financial crisis in 2008, the euro area sovereign debt crisis in 2011, the Covid-19 pandemic in 2020, and the Russian invasion of Ukraine in 2022 (Figure 5). The decomposition of the benchmark SPOT indicator into trigger sources allows for further insights into underlying drivers and helps form a risk narrative. For instance, financial market triggers dominated before and during the global financial crisis and the euro area sovereign debt crisis, while other exogenous triggers increased in importance during the COVID-19 pandemic (Figure 5). More recently, geopolitical triggers spiked following Russia's invasion of Ukraine in 2022 and again during 2025 amidst heightened geopolitical and trade tensions. Overall, these patterns suggest that the proposed AI-based SPOT indicator can measure potential trigger events in a systematic way and provide meaningful information about underlying drivers.

Figure 4: Evolution of aggregate SPOT indicators over time

(a) Frequency of trigger articles

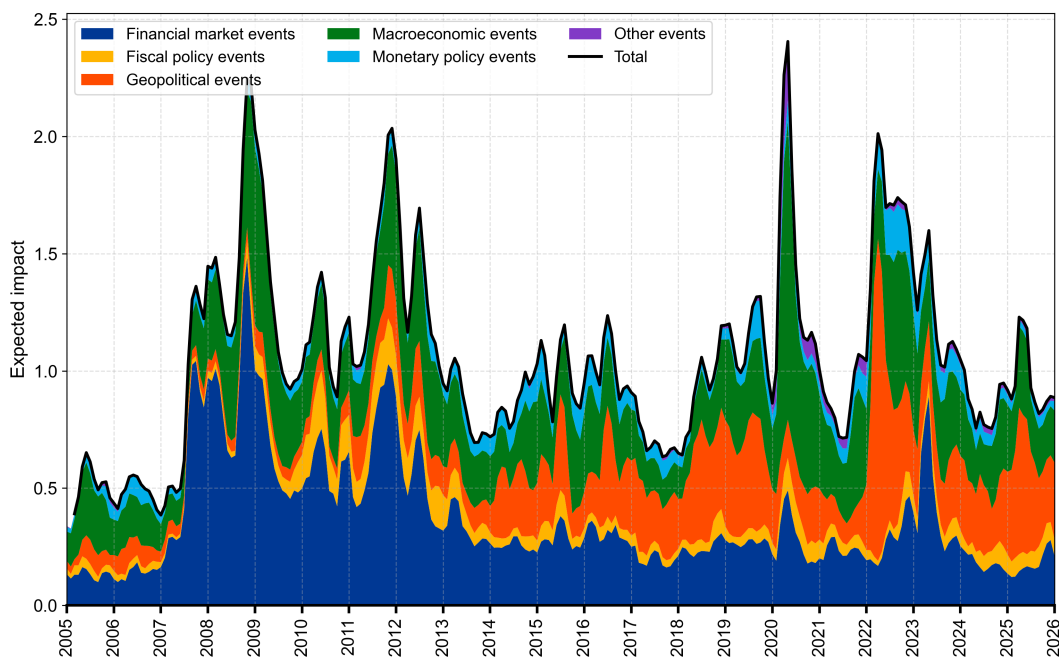


(b) Average probability and average severity of trigger articles



Notes: The trigger frequency indicator shows the share of Class1 articles over the sum of Class0 and Class1 articles. The average probability and severity measures are shown here as averages across only Class1 articles.

Figure 5: Benchmark SPOT indicator with decomposition by trigger source over time

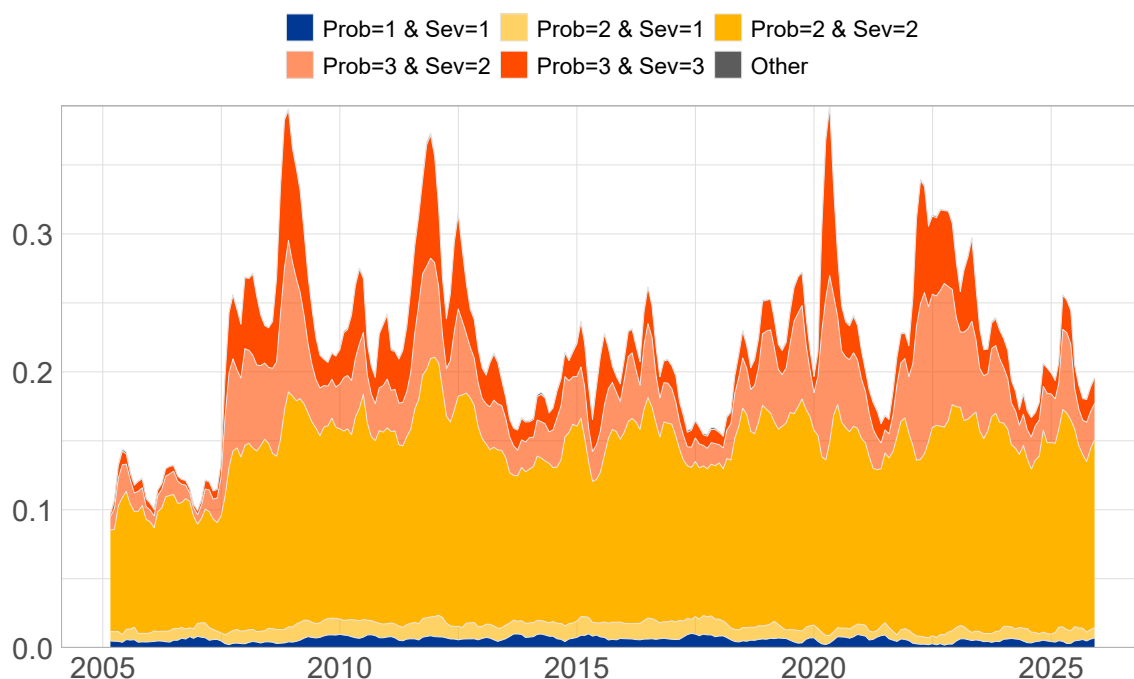


Notes: The benchmark SPOT indicator is calculated as the average expected impact (probability*severity) of trigger events across all articles at a given point in time.

The decomposition of the trigger frequency indicator by the severity and probability of the potential trigger event allows for further insights (Figure 6). Across the full sample, events with a "medium probability" and a "medium severity" account for the largest share of classified trigger articles, reflecting a persistent baseline of concerns about potential trigger events in news coverage.² During major stress episodes like the global financial crisis, the euro area sovereign debt crisis, the COVID-19 pandemic, and the Russian invasion of Ukraine, the composition of classified trigger articles shifts considerably: the share of trigger articles with a high severity rating and/or a high probability rating jumps up, indicating that news narratives intensify. This pattern suggests that monitoring the probability–severity distribution of triggers can capture variation in perceived risk beyond what the trigger frequency alone would convey.

² Such trigger articles are assigned a probability = 2 and a severity = 2.

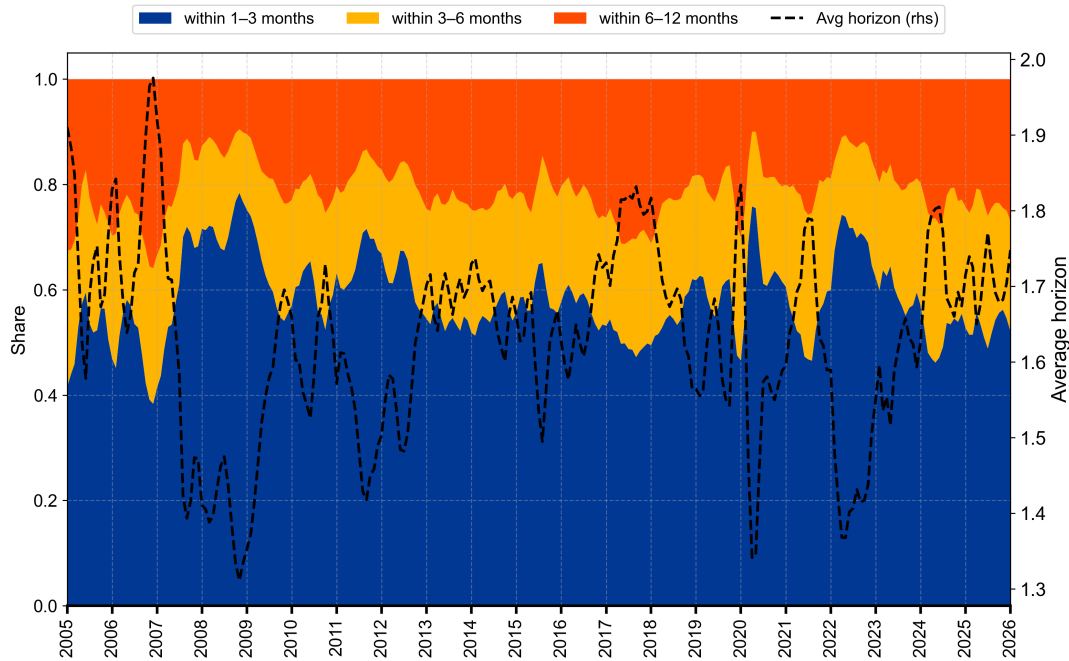
Figure 6: Decomposition of trigger frequency indicator by probability and severity



Notes: The chart shows the frequency of trigger articles by probability and severity classes. Probability is measured on a four-point scale from 0 (unlikely, effectively zero probability) to 3 (high probability). Severity is measured on a four-point scale from 0 (negligible impact) to 3 (very severe or systemic impact).

Beyond the frequency, probability and severity of potential trigger events, the SPOT indicator can also capture a time-horizon dimension. The decomposition of identified trigger articles in [Figure 7](#) shows that trigger-related news coverage is concentrated at short horizons. More than half of trigger articles (58%) are assigned to the short-term bucket (one to three months), with the remainder split roughly equally between the medium-term (three to six months) and longer-term buckets (six to twelve months), as also shown in [Table 1](#). However, the composition of trigger horizons varies markedly over time. The weighted average horizon of trigger articles tends to shorten ahead of and during major stress episodes, including the global financial crisis, the euro area sovereign debt crisis, the COVID-19 pandemic and Russia’s invasion of Ukraine. This is consistent with the notion that, as trigger events become more imminent, media coverage places greater emphasis on near-term adverse developments. From a surveillance perspective, the SPOT horizon provides a complementary signal on the timing of potential triggers, alongside SPOT signals about their perceived probability and severity.

Figure 7: Share of trigger articles by time horizon of potential impact



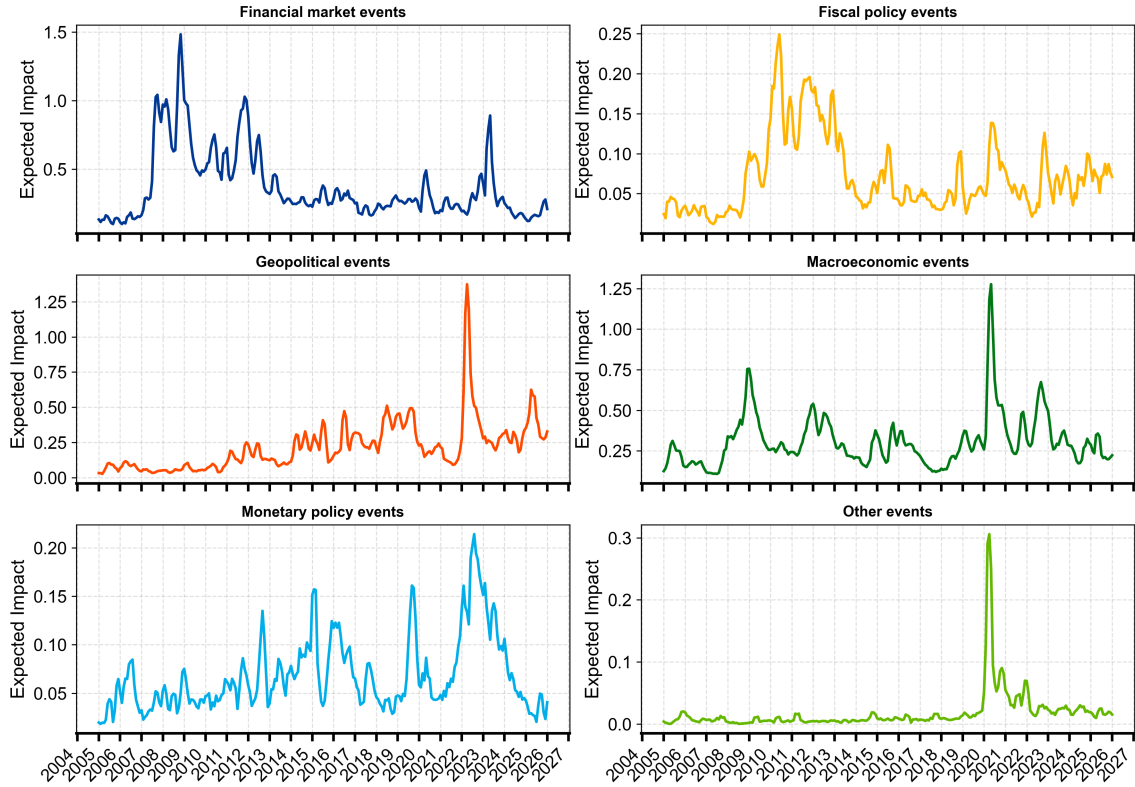
Notes: Horizon = 0 (already occurring events) was dropped from the chart for readability.

4.2 Granular SPOT indicators

Building on the aggregate SPOT decompositions, this subsection further zooms in on four granular dimensions: (i) source-specific SPOT indicators, (ii) higher-frequency SPOT indicators, (iii) horizon-specific SPOT indicators, and (iv) SPOT indicators conditional on probability and severity thresholds.

While the decomposition of the benchmark SPOT indicator into trigger sources can help identify drivers and form an overall risk narrative, sectoral SPOT indicators for individual trigger sources make the underlying drivers even more visible and can be used to complement the benchmark SPOT indicator for monitoring specific trigger sources (Figure 8). For example, the sectoral SPOT indicator for fiscal triggers shows clear peaks in 2010 and 2011, which are less visible in the aggregate decomposition chart (Figure 5). The important role of other exogenous triggers during the Covid-19 pandemic is also much more visible when looking at the sectoral SPOT indicators. The sectoral SPOT indicator for geopolitical triggers also helps to identify a clear upward trend over the past 15 years in the importance of such trigger events (Figure 8). Hence, the sectoral SPOT indicators can provide useful complementary information to the benchmark SPOT indicator.

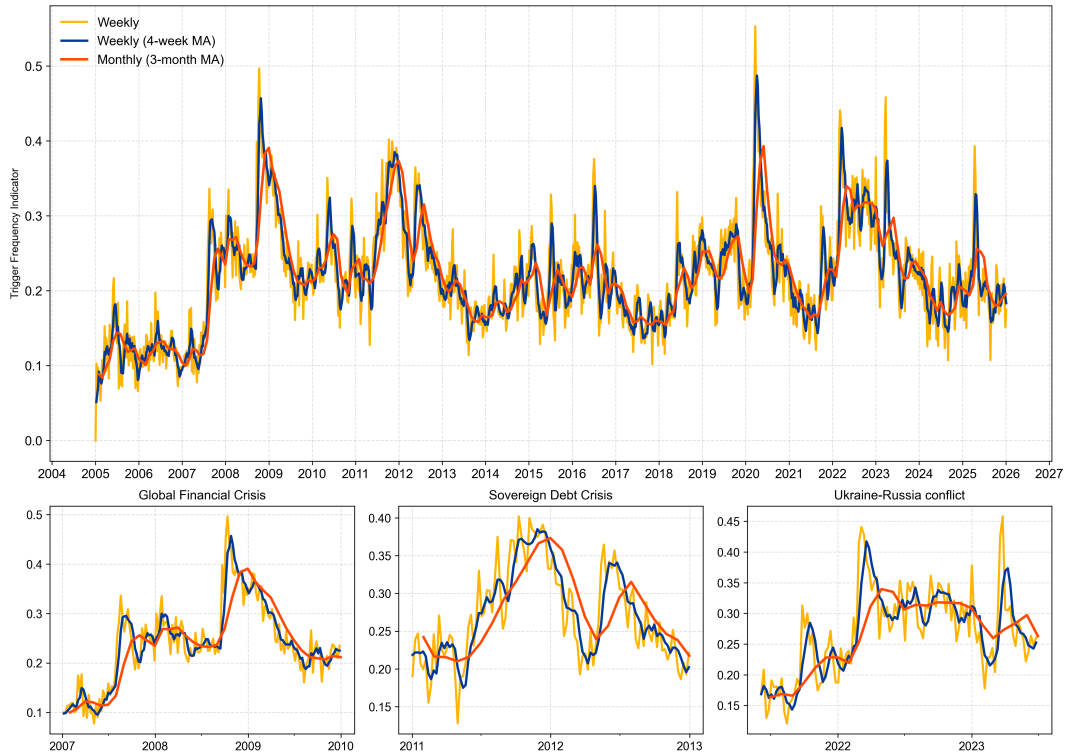
Figure 8: Expected impact indicators for different trigger sources



Notes: The expected impact is calculated across both classes as explained in [subsection 3.3](#).

SPOT indicators can also be constructed at different time aggregations. While the benchmark indicator is calculated as a 3-month moving average of monthly values, higher-frequency measures can provide more timely signals when new triggers emerge. Weekly SPOT indicators capture short-term fluctuations in news coverage of triggers and can signal emerging financial stability risks earlier, whereas monthly and quarterly aggregations smooth temporary movements and highlight more persistent developments in potential trigger events. For example, at the onset of the global financial crisis, the euro area sovereign debt crisis, and Russia’s invasion of Ukraine, higher-frequency SPOT indicators increased earlier than lower-frequency time aggregations, as shown in [Figure 9](#). However, these more timely signals come at the cost of increased volatility, implying a trade-off between responsiveness and noise. Monitoring different time aggregations of SPOT can therefore be useful.

Figure 9: Trigger frequency indicators for different time aggregations

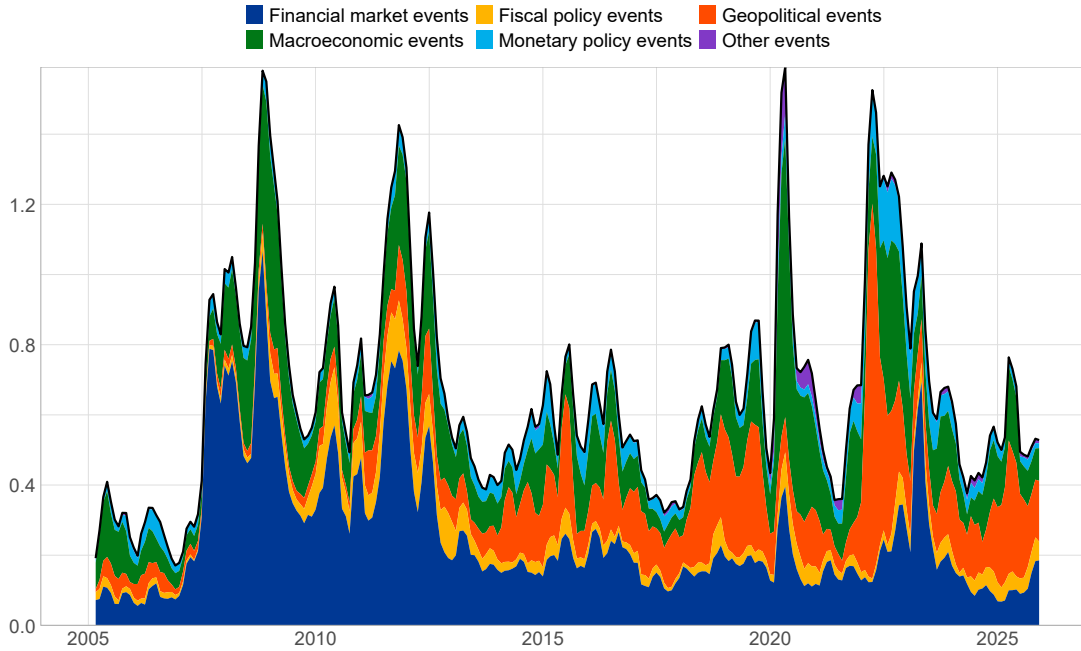


Notes: The trigger frequency indicator is calculated across both classes as explained in [subsection 3.3](#).

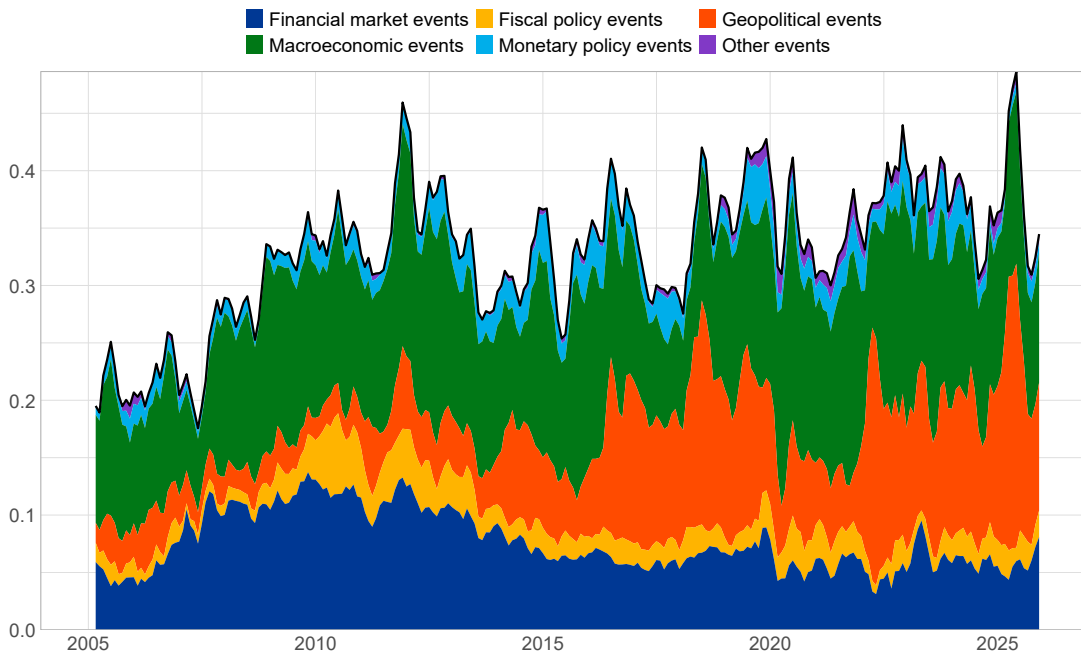
Constructing separate expected impact indicators for short horizons (one to three months) and longer horizons (beyond three months) can further help to identify the immediacy of potential triggers. As shown in [Figure 10](#), the expected impact of potential short-term triggers varies considerably more over time than the expected impact of potential medium-term triggers. In particular, the expected impact of short-term triggers increases significantly around major past crisis episodes, with substantial time-variation in the underlying trigger sources ([Figure 10](#), panel a): during the first half of the sample period, financial market events dominate, with some role also for macroeconomic and fiscal policy trigger events, whereas geopolitical events gain in importance during the second half of the sample, with some role also for macroeconomic and monetary policy events. In contrast, the expected impact of medium-term triggers shows a more secular upward trend over the past 20 years, driven by geopolitical events and macroeconomic events, with peaks around the euro area sovereign debt crisis and the "Liberation day" tariff announcements in 2025 ([Figure 10](#), panel b).

Figure 10: Expected impact indicators for different time horizons

(a) Average expected impact for short horizons (1-3 months)



(b) Average expected impact for longer horizons (more than 3 months)

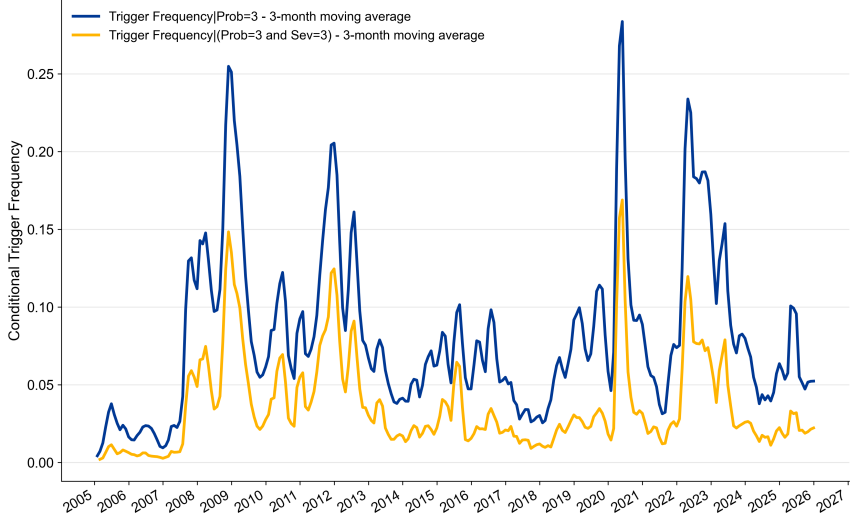


Notes: The indicators are calculated as the average expected impact (probability*severity) of trigger events for the respective horizon group across all articles at a given point in time.

Finally, granular SPOT indicators conditioning on certain probability and severity thresholds can also provide further insights. For example, the trigger frequency conditioning on

high-probability events is consistently larger than the frequency conditioning also on high-severity events, suggesting that news coverage more frequently signals an increased probability of trigger events materialising rather than also having a highly severe macro-financial impact (Figure 11). In addition, the frequency of trigger events classified as highly probable and highly severe mainly shoots up during episodes of major stress, whereas the frequency of high-probability trigger events increases more gradually and often rises somewhat in advance. These patterns suggest that increases in the frequency of triggers with a high-probability and high-severity signals more imminent and consequential threats, while increases in the frequency of high-probability events alone can signal earlier shifts in potential trigger events. Moreover, the frequency of high-probability and high-severity triggers appears less responsive to news with uncertain outcomes, as illustrated by the relatively muted response to the "Liberation day" tariff announcements in 2025, which were widely discussed but not uniformly assessed as having a very severe macro-financial impact.

Figure 11: Trigger frequency indicators for different probability and severity thresholds



Notes: The conditional trigger frequency indicator shows the share of Class1 articles given the specific conditions hold true over the sum of Class0 and conditioned Class1 articles (see subsection 3.3).

5 Evaluation of SPOT indicators

5.1 Comparison to other risk indicators from the literature

In this section, we compare SPOT indicators to selected indicators proposed in the literature for measuring financial stress, uncertainty, geopolitical risk or other relevant aspects that can relate to trigger events. The indicators are listed in Table 3 together with selected summary statistics

and Figure 12 presents pairwise correlations among the various indicators and annual real GDP growth for the euro area. A number of observations are worth noting.

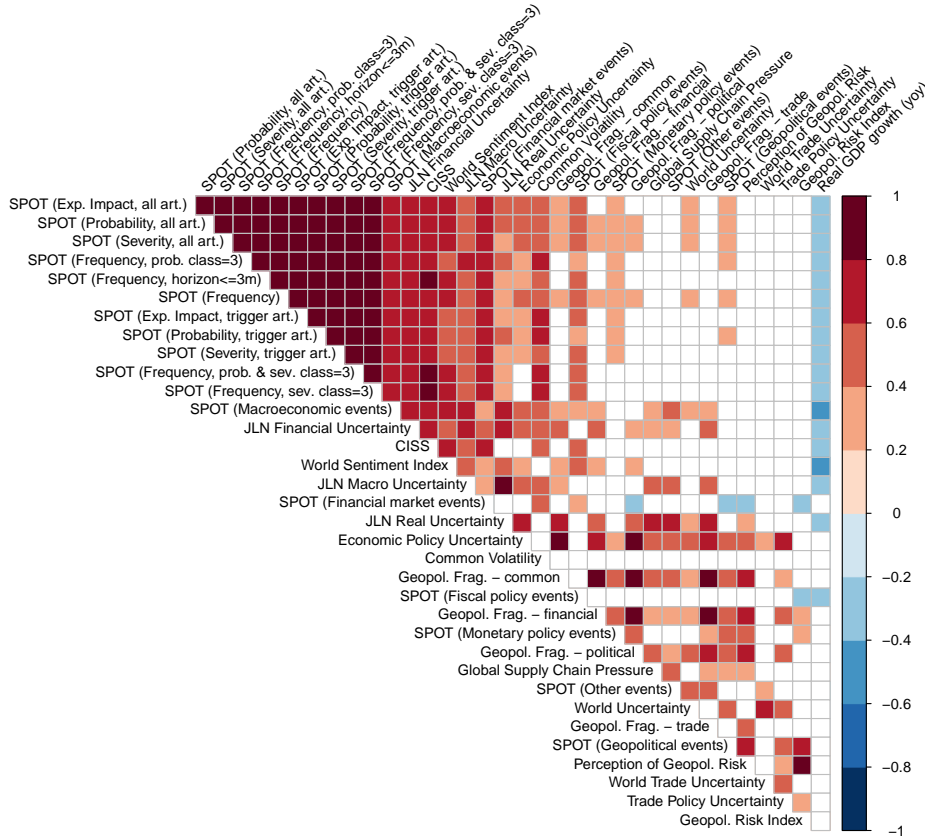
Table 3: Overview of indicators evaluated in the empirical exercise

| Indicator | Source | Mean | St. deviation | Median | 25th perc. | 75th perc. | Min | Max | Latest obs. | |
|--|------------------------------------|---------------------------------|---------------|--------|------------|------------|--------|--------|-------------|---------|
| SPOT (Frequency) | This paper | 0.22 | 0.06 | 0.22 | 0.18 | 0.24 | 0.10 | 0.39 | 2025 Q4 | |
| SPOT (Probability, all art.) | | 0.50 | 0.16 | 0.48 | 0.40 | 0.57 | 0.20 | 0.98 | 2025 Q4 | |
| SPOT (Probability, trigger art.) | | 2.25 | 0.12 | 2.24 | 2.18 | 2.30 | 2.02 | 2.52 | 2025 Q4 | |
| SPOT (Severity, all art.) | | 0.45 | 0.14 | 0.43 | 0.36 | 0.51 | 0.19 | 0.86 | 2025 Q4 | |
| SPOT (Severity, trigger art.) | | 2.04 | 0.08 | 2.02 | 1.99 | 2.09 | 1.89 | 2.21 | 2025 Q4 | |
| SPOT (Exp. Impact, all art.) | | 1.06 | 0.39 | 1.00 | 0.82 | 1.22 | 0.39 | 2.22 | 2025 Q4 | |
| SPOT (Exp. Impact, trigger art.) | | 4.72 | 0.42 | 4.65 | 4.45 | 4.98 | 3.90 | 5.71 | 2025 Q4 | |
| SPOT (Frequency, horizon<=3m) | | 0.13 | 0.05 | 0.12 | 0.10 | 0.16 | 0.04 | 0.31 | 2025 Q4 | |
| SPOT (Frequency, prob. class=3) | | 0.07 | 0.04 | 0.06 | 0.04 | 0.08 | 0.01 | 0.20 | 2025 Q4 | |
| SPOT (Frequency, sev. class=3) | | 0.03 | 0.02 | 0.02 | 0.02 | 0.04 | 0.00 | 0.10 | 2025 Q4 | |
| SPOT (Frequency, prob. & sev. class=3) | | 0.03 | 0.02 | 0.02 | 0.01 | 0.04 | 0.00 | 0.09 | 2025 Q4 | |
| SPOT (Financial market events) | | 0.37 | 0.25 | 0.28 | 0.22 | 0.47 | 0.10 | 1.24 | 2025 Q4 | |
| SPOT (Fiscal policy events) | | 0.07 | 0.05 | 0.06 | 0.04 | 0.09 | 0.01 | 0.25 | 2025 Q4 | |
| SPOT (Geopolitical events) | | 0.22 | 0.18 | 0.19 | 0.09 | 0.29 | 0.03 | 1.16 | 2025 Q4 | |
| SPOT (Macroeconomic events) | | 0.31 | 0.15 | 0.28 | 0.21 | 0.37 | 0.11 | 0.99 | 2025 Q4 | |
| SPOT (Monetary policy events) | | 0.07 | 0.04 | 0.06 | 0.04 | 0.08 | 0.02 | 0.19 | 2025 Q4 | |
| SPOT (Other events) | | 0.02 | 0.03 | 0.01 | 0.01 | 0.02 | 0.00 | 0.29 | 2025 Q4 | |
| CISS | | Kremer et al. (2012) | 0.17 | 0.19 | 0.08 | 0.03 | 0.27 | 0.00 | 0.90 | 2025 Q4 |
| Common Volatility | | Engle and Campos-Martins (2023) | 0.57 | 0.16 | 0.55 | 0.47 | 0.66 | 0.19 | 1.10 | 2025 Q2 |
| Economic Policy Uncertainty | | Baker et al. (2016) | 162.85 | 72.65 | 151.87 | 111.72 | 210.02 | 55.39 | 370.51 | 2025 Q1 |
| Geopol. Frag. - financial | Fernandez-Villaverde et al. (2024) | -0.63 | 0.63 | -0.78 | -0.95 | -0.23 | -1.68 | 0.99 | 2024 Q1 | |
| Geopol. Frag. - political | Fernandez-Villaverde et al. (2024) | 1.51 | 1.30 | 1.72 | 0.45 | 2.52 | -0.66 | 3.48 | 2024 Q1 | |
| Geopol. Frag. - trade | Fernandez-Villaverde et al. (2024) | -0.75 | 0.33 | -0.89 | -0.97 | -0.68 | -1.12 | 0.28 | 2024 Q1 | |
| Geopol. Risk Index | Caldara and Iacoviello (2022) | 100.25 | 26.04 | 90.85 | 84.92 | 109.19 | 69.69 | 224.60 | 2025 Q2 | |
| Geopol. Frag. - common | Fernandez-Villaverde et al. (2024) | -0.32 | 0.45 | -0.44 | -0.52 | -0.14 | -0.96 | 0.76 | 2024 Q1 | |
| Global Supply Chain Pressure | FED New York | 0.16 | 1.06 | -0.12 | -0.48 | 0.33 | -1.32 | 4.26 | 2025 Q2 | |
| JLN Financial Uncertainty | Ludvigson et al. (2021) | 0.99 | 0.06 | 0.99 | 0.94 | 1.03 | 0.90 | 1.13 | 2024 Q3 | |
| JLN Macro Uncertainty | Ludvigson et al. (2021) | 0.93 | 0.07 | 0.91 | 0.88 | 0.97 | 0.84 | 1.13 | 2024 Q3 | |
| JLN Real Uncertainty | Ludvigson et al. (2021) | 0.89 | 0.05 | 0.86 | 0.86 | 0.91 | 0.83 | 1.10 | 2024 Q3 | |
| Perception of Geopol. Risk | Alonso-Alvarez et al. (2025) | 99.83 | 48.44 | 92.14 | 66.69 | 105.54 | 55.02 | 358.89 | 2024 Q4 | |
| Trade Policy Uncertainty | Caldara et al. (2019) | 70.25 | 104.00 | 34.60 | 26.64 | 61.87 | 20.78 | 782.33 | 2025 Q2 | |
| World Sentiment Index | Ahir et al. (2022) | 0.82 | 0.68 | 0.82 | 0.49 | 1.30 | -1.40 | 2.16 | 2025 Q3 | |
| World Trade Uncertainty | Ahir et al. (2022) | 11.04 | 29.53 | 1.03 | 0.15 | 3.73 | 0.00 | 174.34 | 2025 Q1 | |
| World Uncertainty | Ahir et al. (2022) | 22.06 | 9.31 | 19.65 | 15.10 | 26.46 | 8.24 | 174.34 | 2025 Q1 | |

Notes: Summary statistics are computed on quarterly observations starting from 2005 Q1. Indicators available at monthly frequency are transformed to quarterly values by computing the moving average over 3 months. Sectoral SPOT indicators reflect the expected impact for that trigger source. World Uncertainty is in thousands.

First, SPOT indicators are negatively correlated with real GDP growth, which is consistent with the notion that an elevated trigger probability and severity is associated with weaker economic activity (Figure 12). Second, various SPOT indicators exhibit high correlations among each other. Nonetheless, this is mostly true for the aggregate SPOT indicators, while sectoral SPOT indicators exhibit lower correlations. Third, the benchmark SPOT displays positive but imperfect correlations with established indicators such as the Composite Indicator of Systemic Stress (CISS), the Geopolitical Risk Indicator (GPR), and the Economic Policy Uncertainty Indicator (EPU). This shows that while these measures share common dimensions, SPOT captures information beyond what is reflected in widely used financial stress and risk indicators. Fourth, sectoral SPOT indicators cluster with their thematically relevant counterparts — SPOT for geopolitical events groups with GPR, trade policy uncertainty, and geopolitical fragmentation measures, while SPOT for financial market events groups with Common Volatility, CISS and JLN financial uncertainty. This pattern lends support to the validity of the LLM-based classification of trigger sources.

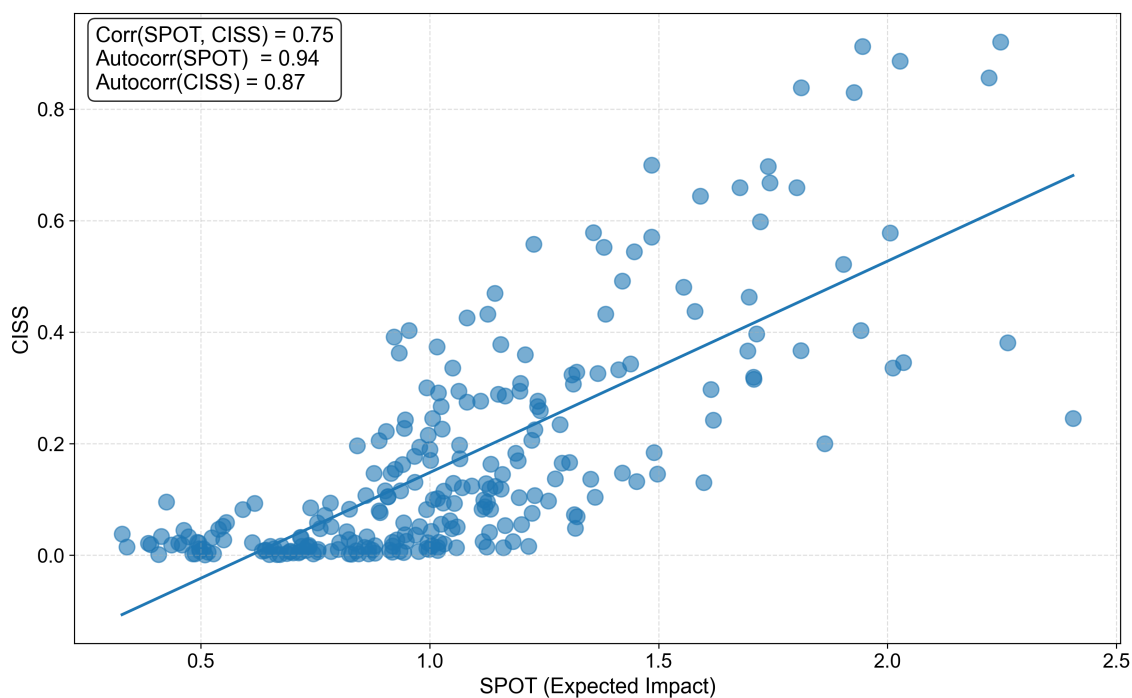
Figure 12: Correlations between SPOT indicators and other indicators



Notes: Computed on quarterly sample 2005Q1-2024Q4 using 3 month averages of monthly values. White squares denote correlations that are not significant at the 99% level. For the World Sentiment Index, negative values are used so that increases denote higher risk. Indicators are ordered based on the first principal component loadings.

Compared to the CISS (Kremer et al., 2012), a widely used financial stress indicator for the euro area, the benchmark SPOT indicator exhibits a clear positive correlation of 0.75, yet the relationship is imperfect (Figure 13). This suggests that the SPOT indicator provides complementary news-based information beyond what is captured by financial market conditions. Notably, the SPOT indicator displays greater persistence than the CISS, consistent with its narrative nature: while the CISS, which is constructed from volatility measures and cross-correlations across financial market segments, tends to respond swiftly to new information, the SPOT indicator reflects the more gradual build-up and unwinding of trigger-related narratives in financial news. Moreover, a key advantage of the SPOT indicator over purely market-based stress measures is its better interpretability. Whereas the CISS captures how financial markets react in a given situation, through elevated volatility and tighter cross-market co-movements, it does not reveal the underlying narratives for the observed stress. In contrast, the SPOT indicator’s decomposition into trigger sources provides a structured narrative of what is driving the increase in potential trigger events, adding information for financial stability surveillance purposes.

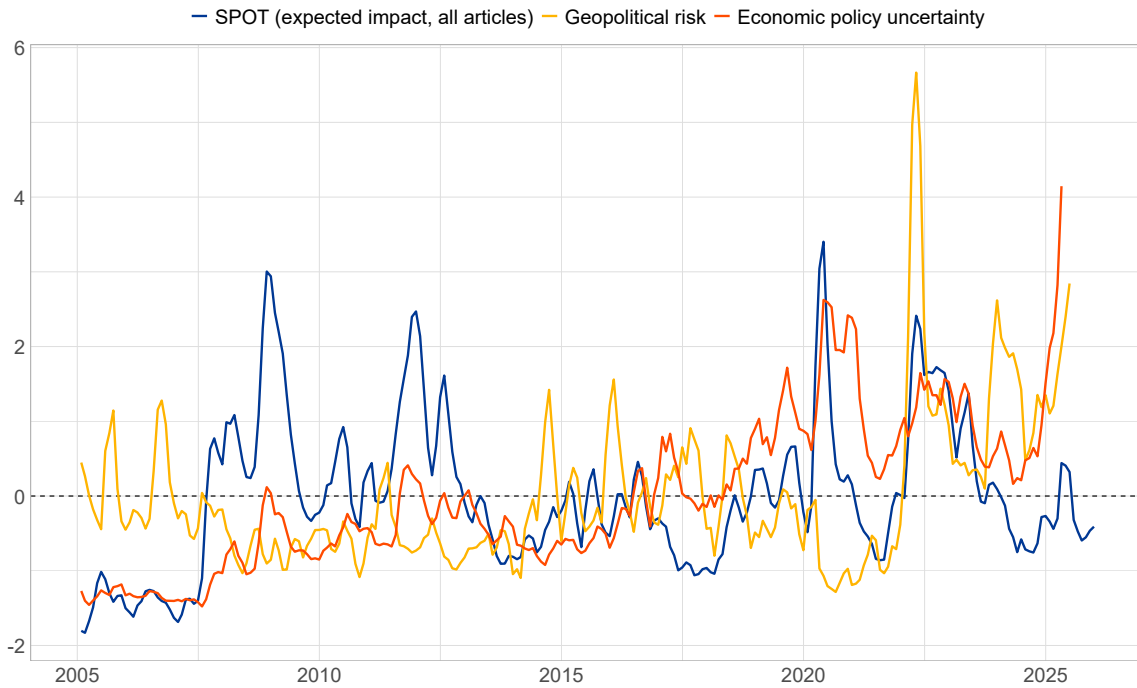
Figure 13: Visual comparison of the benchmark SPOT indicator with the CISS



Notes: The benchmark SPOT (Expected Impact) and CISS values are monthly observations.

We also compare the benchmark SPOT indicator to two prominent text-based indicators: the GPR of [Caldara and Iacoviello \(2022\)](#) and the EPU of [Baker et al. \(2016\)](#). The comparison reveals episodes of clear co-movement ([Figure 14](#)). For instance, both the benchmark SPOT indicator and EPU rise markedly during the COVID-19 pandemic in 2020 and amid heightened trade policy tensions in 2025, while the SPOT indicator and the GPR index increase together following Russia's invasion of Ukraine in 2022. However, there are also notable episodes where the SPOT indicator captures potential triggers that these more specialised indicators do not. For example, during the global financial crisis and the euro area sovereign debt crisis the SPOT indicator rises substantially, whereas neither the GPR nor the EPU fully reflect the severity of these episodes, given that they were not primarily driven by geopolitical events or policy uncertainty. In addition, even in periods of co-movement, the amplitudes of the respective indicators can differ, reflecting their distinct underlying construction. One useful feature of the SPOT indicator is that it classifies a trigger article only when it suggests a potentially meaningful adverse effect on the economy or financial system over a forward-looking horizon, rather than simply flagging geopolitical or policy-uncertainty language in news. By combining trigger event identification with an assessment of the likely macro-financial impact, SPOT can provide information that complements standard dictionary-based risk indicators.

Figure 14: Visual comparison of the benchmark SPOT indicator with the GPR and the EPU

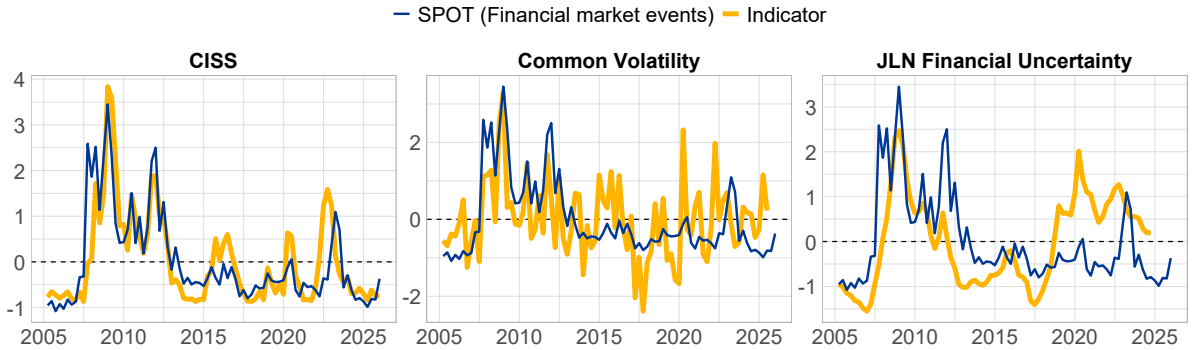


Notes: Indicators are standardised to have a mean of zero and a standard deviation of one.

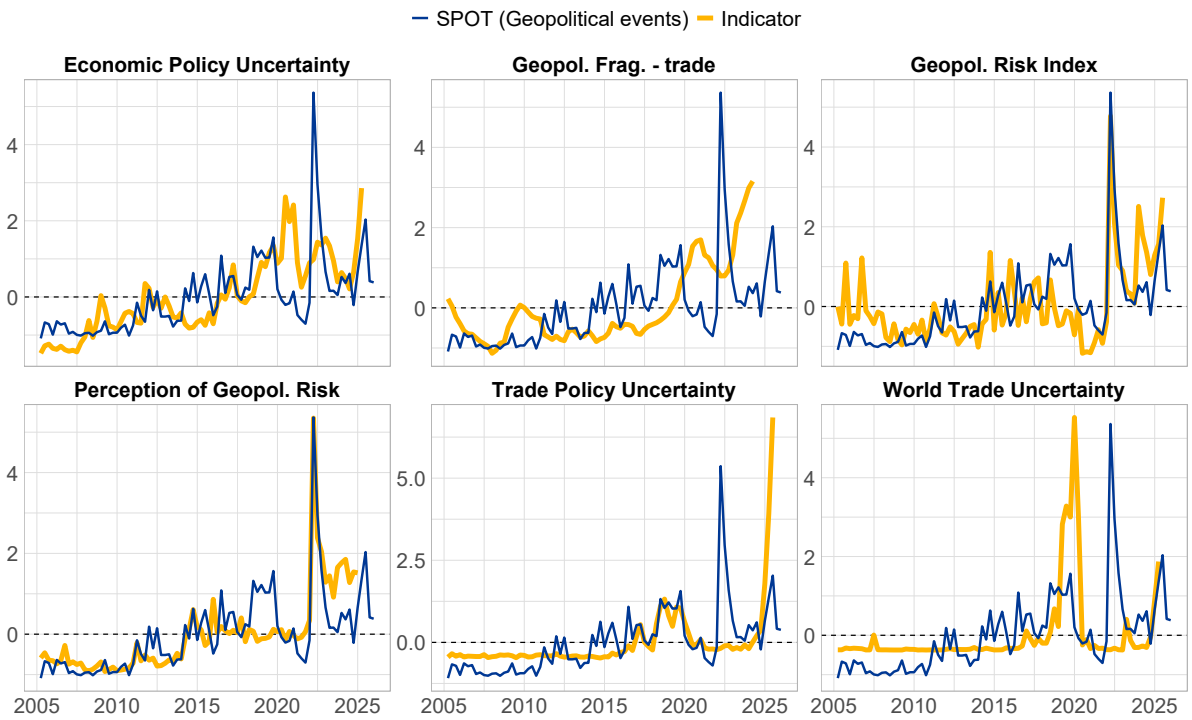
Finally, several sectoral SPOT indicators are compared to established risk and uncertainty indicators (Figure 15). Three observations are worth noting. First, consistent with the correlation analysis above, the sectoral SPOT indicators broadly co-move with their thematically closest counterpart indicators: for example with the CISS, common volatility, and JLN financial uncertainty in the case of the financial market SPOT; with the GPR, EPU, and trade policy uncertainty for the geopolitical SPOT; and with the JLN macro and real uncertainty for the macroeconomic events SPOT. This thematic co-movement corroborates the consistency of the LLM-based trigger source classification. Second, the co-movement is loose rather than exact. There are episodes during which the sectoral SPOT indicators rise while the corresponding other indicators remain comparatively muted, and vice versa. Such occasional divergences suggest that the LLM-based SPOT indicators may capture aspects of trigger perceptions that may not be fully reflected in market-based or dictionary-based indicators. Third, despite these differences, the overall time-series profile of the sectoral SPOT indicators is broadly consistent with that of established measures, indicating that the SPOT indicators do not generate spurious signals unrelated to the risk environment.

Figure 15: Visual comparison of sectoral SPOT indicators thematically aligned indicators

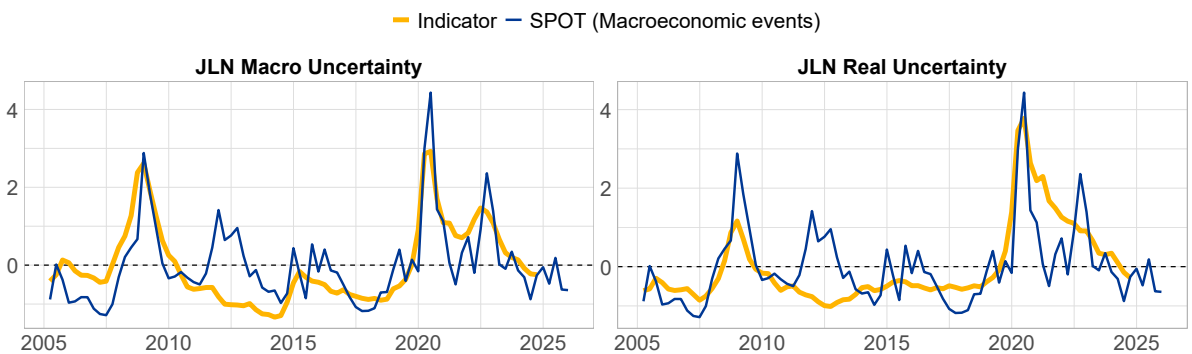
(a) Financial market events



(b) Geopolitical events



(c) Macroeconomic events



Notes: Indicators are standardised to have a mean of zero and a standard deviation of one. Sectoral SPOT indicators shown capture expected impact of all articles.

5.2 SPOT performance evaluation in growth-at-risk models

Design of the growth-at-risk exercises

To assess the information content of SPOT indicators compared to other risk indicators more formally, we employ a growth-at-risk model framework that can be used to estimate future downside risks to the economy (Adrian et al., 2022, 2019; Lang et al., 2025). In particular, we estimate quarterly panel quantile local projection models for various horizons h as follows:

$$Q_{y_{i,t+h},\tau} = \alpha_i^{h,\tau} + \rho^{h,\tau} y_{i,t} + \beta^{h,\tau} \mathbf{X}'_{i,t} + \varepsilon_{i,t+h,\tau},$$

where quantiles are denoted by τ , $y_{i,t+h}$ is the average real GDP growth rate in country i between time period t and $t + h$, $\mathbf{X}'_{i,t}$ is a vector of explanatory variables, while $\alpha_i^{h,\tau}$ are country fixed effects, and $\varepsilon_{i,t+h,\tau}$ denotes an error term. We apply a two-step estimation approach for panel quantile regressions following Canay (2011). In the first step, we estimate the unobserved fixed effects using a within estimator, assuming that country fixed effects are constant across quantiles. In the second step, we run a standard conditional quantile regression on the dependent variable that has been adjusted by subtracting the fixed-effect estimates from the first step (see Canay, 2011, for further details).

To compare the performance of different indicators, we use the tick loss as our measure of model fit, a criterion widely used to assess the accuracy of value-at-risk models. The majority of growth-at-risk studies that adopt a formal evaluation metric make use of the tick loss function (Brownlees and Souza, 2021; Carriero et al., 2020; Figueres and Jarocinski, 2020; Giacomini and Komunjer, 2005). The tick loss is the value of the objective function that is minimized by the quantile regression. More precisely, it is computed as follows:

$$TL_{h,\tau} = (\tau - \mathbf{1}(\hat{\varepsilon}_{h,\tau} < 0))\hat{\varepsilon}_{h,\tau},$$

where $\mathbf{1}()$ denotes the indicator function, τ is the quantile of interest, and $\hat{\varepsilon}_{h,\tau}$ is the residual (i.e., actual value minus predicted value).

As a benchmark model, we estimate a panel quantile regression with current annual and quarterly real GDP growth as explanatory variables as well as the systemic risk indicator (SRI) of Lang et al. (2019), its lag, and an interaction term for when the SRI is positive to capture financial stability vulnerabilities. As shown by Lang et al. (2025), the SRI outperforms other vulnerability indicators in terms of in-sample explanatory power and out-of-sample forecasting performance for medium-term growth-at-risk in euro area countries and therefore constitutes a

key benchmark against which the additional information content of other indicators should be evaluated. To account for the exogenous nature of the covid pandemic shock, we add dummies for the following periods: 2020 Q1, ..., 2020 Q4, and 2021 Q2. The models are estimated on a sample of euro area countries³ over the period 2005 Q1 to 2024 Q1.

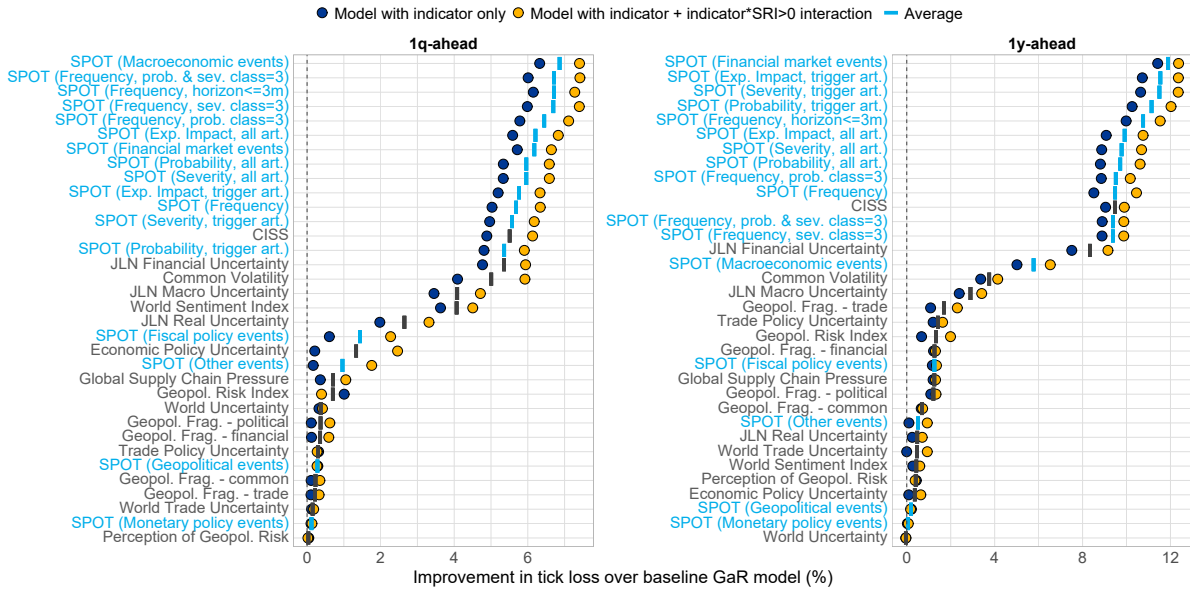
We focus on the lower tail of the future real GDP growth distribution by estimating the 10th quantile ($\tau = 0.1$) for 1-quarter and 1-year ahead horizons ($h = 1, 4$). Each potential trigger indicator is added to the benchmark model one at a time. To make the results comparable, we balance observations across all models (i.e. the dataset is balanced within each country). We also consider specifications that interact the indicator with the *SRI* whenever the *SRI* is positive, to capture the potential amplification of triggers when vulnerabilities are elevated. For each specification, we compute the percentage improvement in the average tick loss relative to the benchmark model.

Performance of various indicators in growth-at-risk exercises

Key results from the growth-at-risk exercises are reported in [Figure 16](#). Overall, SPOT indicators significantly improve the performance of the benchmark model and consistently outperform other standard financial stress and risk indicators. This holds for both the 1-quarter ahead and 1-year ahead horizon, with improvements compared to the benchmark model fit of more than 6% and 10% respectively. Given the high correlation between various aggregate SPOT indicators documented in [Figure 12](#), it is not surprising that their growth-at-risk model performance is similar. Nevertheless, SPOT indicators that combine information about the frequency of triggers with information about the probability/severity of triggers tend to slightly dominate SPOT indicators that are solely based on the trigger frequency. In addition, SPOT indicators that combine information about the probability with information about the severity (expected impact) perform slightly better than indicators solely based on the probability or the severity. These findings indicate that combining different trigger dimensions in aggregate SPOT indicators seems desirable. The benchmark SPOT indicator (average expected impact) is among the top six best-performing SPOT indicators for both growth-at-risk projection horizons.

³ In addition to euro area countries, we include DK and SE but exclude HR and BG due to data availability.

Figure 16: Performance of SPOT and other indicators in growth-at-risk models



Notes: Improvements in tick loss function are in percentage relative to the benchmark model with GDP and the SRI (including its lag and a dummy variable for cases when $SRI > 0$). Indicators are ordered by largest average improvement in model performance (i.e. percentage decrease in the tick loss relative to the benchmark model) for the two sets of model specifications (indicator only and indicator*SRI>0 interaction models).

The best-performing SPOT indicator for the 1-quarter ahead horizon is the expected impact of macroeconomic trigger events, while for the 1-year ahead horizon the best-performing SPOT indicator is the expected impact of financial market trigger events (Figure 16). However, neither of these sectoral SPOT indicators outperform the benchmark SPOT indicator at both horizons. Other granular SPOT indicators that slightly outperform the benchmark SPOT indicator for the 1-quarter ahead horizon are the frequency of triggers with a high probability and/or high severity. However, at the 1-year ahead horizon this slight out-performance vanishes. Another granular SPOT indicator that performs well at both growth-at-risk horizons is the frequency of immediate trigger events (trigger horizon ≤ 3 months). Among other established risk indicators, those related to financial markets provide the largest improvements in growth-at-risk model fit (CISS, JLN Financial Uncertainty, and Common Volatility). However, their performance remains below the performance of the benchmark SPOT. Most other risk indicators are significantly worse than the benchmark SPOT. The strong growth-at-risk model improvement of financial market indicators (including the sectoral SPOT indicator for financial market events) is likely related to the fact that the estimation sample includes two crises driven by financial markets (the global financial crisis and the sovereign debt crisis).

As shown in Figure 16, the growth-at-risk model performance of SPOT indicators improves further, by around 2 percentage points, when interaction terms between the SPOT indicators

and the SRI are included (see yellow vs blue dots). However, the overall growth-at-risk performance ordering of various SPOT indicators does not change much when allowing for such interaction terms. These results indicate that there are important interactions and amplification mechanisms between underlying financial vulnerabilities, as measured by the SRI, and potential trigger events captured by SPOT: while triggers act as immediate catalysts for stress, the differential impact on the lower tail of the growth distribution is driven by pre-existing macro-financial vulnerabilities. This highlights the importance of a holistic approach to monitoring financial stability risks, combining information about vulnerabilities with information about potential trigger events.

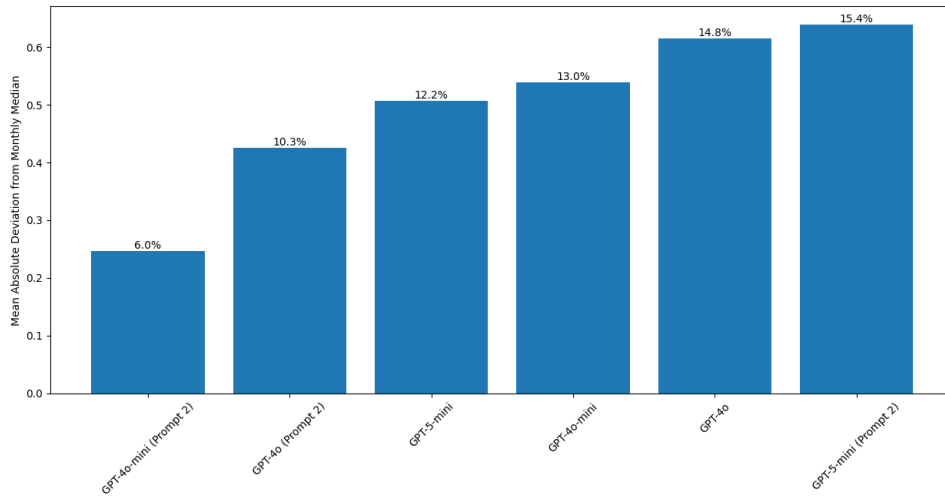
The findings regarding the growth-at-risk performance of various SPOT indicators are robust to changing the sample period or looking at different quantiles, as shown in Annex [A.2](#).

5.3 Robustness and sensitivity analysis of SPOT

In this section we assess the sensitivity of SPOT results to changes in the prompt design and the LLM version used. For this purpose we randomly draw 100 articles per month from the pre-filtered dataset of articles used in stage 2 of our prompting pipeline. We evaluate these articles using GPT-4o-mini, GPT-4o, and GPT-5-mini using two distinct prompt specifications.

To quantify cross-model variation, we first construct a "consensus" SPOT series by taking the monthly median value of the SPOT indicator across all model and prompt specifications. This "consensus" SPOT series serves as a benchmark representing the central tendency across models and prompts. For each model-prompt specification, we then compute the average absolute %-deviation of the model-prompt-specific SPOT indicator from this "consensus" SPOT indicator over all months in the sample. As shown in [Figure 17](#), there is only moderate variation in the SPOT indicator across model-prompt specifications. Relative deviations range from 6.0% for GPT-4o-mini (Prompt 2) to 15.4% for GPT-5-mini (Prompt 2), with most models deviating from the "consensus" SPOT by approximately 10–15% on average. These magnitudes suggest that, while article assessments are somewhat model- and prompt-dependent, the overall extracted signals remain broadly consistent across model-prompt specifications. No model appears structurally detached from the central tendency ("consensus"), indicating a stable common component in the generated SPOT indicators despite differences in model architecture and prompting.

Figure 17: Mean absolute deviation of SPOT versions based on different LLMs and prompts



Notes: Deviations are measured relative to the median SPOT value across all model-prompt specifications ("consensus"). All model-prompt combinations are run on a common sample of 100 randomly drawn articles per month.

Another potential robustness concern is whether the LLM implicitly relies on future information when classifying articles and assessing severity, horizon, and probability of the triggers. We partially mitigate this concern by explicitly instructing the model in the system prompt to base its evaluation solely on the information contained in each article. Hence, our framework is designed to classify and structure the content of individual texts rather than to generate forecasts. Consequently, even if some forward-looking attributes were affected by leakage, this would not undermine the broader objective of extracting structured information from news as a complement to traditional indicators. However, we acknowledge that prompting alone cannot fully eliminate the risk of data leakage.

The different knowledge cut-offs of the underlying models provides a potential cross-check for whether data leakage is an issue or not. For GPT-4 and GPT-4o-mini, the knowledge cut-off is October 2023, while for GPT-5-mini it is the end of May 2024. We find that the deviations of the different SPOT indicators from the "consensus" in between these two cut-offs are not systematically different from the deviations observed over the full sample. This indicates that use of outside knowledge by the LLMs to assess articles does not seem to be a problem. As a further robustness exercise, we also conduct experiments in which the publication date of the articles is withheld from the model. The article assessments remain largely unchanged compared to when the LLM is also fed with the publication date of the article, supporting the view that the model does not systematically rely on temporal cues or future information.

6 Conclusion

In this paper we developed SPOT, an AI-based indicator that measures the potential **S**everity and **P**robability **O**f **T**riggers (**SPOT**) for financial stability risks based on news. We showed that our benchmark SPOT indicator spikes ahead of major historical episodes of financial instability, such as the global financial crisis, the euro area sovereign debt crisis, the Covid-19 pandemic, or the Russian invasion of Ukraine. Moreover, the decomposition of SPOT by trigger source allows for further insights into underlying drivers and helps form a risk narrative. Sectoral SPOT indicators for individual trigger sources, e.g. for geopolitical or fiscal triggers, make the underlying drivers even more visible and can be used to complement the benchmark SPOT indicator for monitoring specific trigger sources. We also showed that various SPOT indicators contain useful information about future downside risks to the economy and outperform other commonly used indicators capturing financial stress, geopolitical risk, policy uncertainty or volatility in euro area panel growth-at-risk models. Hence, the SPOT framework provides a comprehensive view of the dynamics, drivers, and characteristics of perceived triggers for financial stability risks.

Overall, the results presented in this paper suggest that AI-based signal extraction from text offers a promising avenue to enhance the monitoring of financial stability risks. Recent research shows that machine learning and LLMs can be used to extract relevant information from unstructured data to construct indicators of sentiment, uncertainty, and macro-financial risks. Building on these advances, the SPOT indicator presented in this paper demonstrates how LLMs can be applied to financial news to provide a systematic and interpretable measure of potential trigger events, thereby complementing existing vulnerability indicators for monitoring financial stability risks. Looking ahead, a promising avenue for applied financial stability research appears to be the application of similar approaches to new or previously underutilised data sources, including financial disclosures, supervisory information, market intelligence findings, or social media posts. While challenges related to model reliability, interpretability, and data governance remain, AI-based approaches have the potential to become an increasingly important component of financial stability surveillance frameworks.

References

- Adrian, Tobias, Federico Grinberg, Nellie Liang, Sheheryar Malik, and Jie Yu**, “The Term Structure of Growth-at-Risk,” *American Economic Journal: Macroeconomics*, July 2022, 14 (3), 283–323.
- , **Nina Boyarchenko, and Domenico Giannone**, “Vulnerable Growth,” *American Economic Review*, April 2019, 109 (4), 1263–1289.
- Ahir, Hites, Nicholas Bloom, and Davide Furceri**, “The World Uncertainty Index,” NBER Working Papers 29763, National Bureau of Economic Research, Inc Feb 2022.
- Alessi, Lucia and Carsten Detken**, “Quasi real time early warning indicators for costly asset price boom/bust cycles: A role for global liquidity,” *European Journal of Political Economy*, September 2011, 27 (3), 520–533.
- Alonso-Alvarez, Irma, Marina Diakonova, and Javier J. Perez**, “Rethinking GPR: The sources of geopolitical risk,” Working Papers 2522, Banco de España May 2025.
- Audrino, Francesco, Jessica Gentner, and Simon Stalder**, “Quantifying Uncertainty: A New Era of Measurement through Large Language Models,” *Swiss Finance Institute Research Paper Series*, July 2024, (24-68).
- Baker, Scott R., Nicholas Bloom, and Steven J. Davis**, “Article Navigation Journal Article Editor’s Choice Measuring Economic Policy Uncertainty,” *The Quarterly Journal of Economics*, November 2016, 131 (4), 1593–1636.
- Bond, Shaun A., Hayden Klok, and Min Zhu**, “Large Language Models and Financial Market Sentiment,” *SSRN Online Library*, October 2023.
- Borio, Claudio and Mathias Drehmann**, “Assessing the risk of banking crises - revisited,” *BIS Quarterly Review*, March 2009, pp. 29–46.
- **and Philip Lowe**, “Asset prices, financial and monetary stability: exploring the nexus,” BIS Working Papers 114, Bank for International Settlements July 2002.
- Brownlees, Christian and André B.M. Souza**, “Backtesting global Growth-at-Risk,” *Journal of Monetary Economics*, 2021, 118 (C), 312–330.

Caldara, Dario and Matteo Iacoviello, “Measuring Geopolitical Risk,” *American Economic Review*, April 2022, 112 (4), 1194–1225.

—, —, **Patrick Molligo, Andrea Prestipino, and Andrea Raffo**, “Does Trade Policy Uncertainty Affect Global Economic Activity?,” FEDS Notes 2019-09-04, Board of Governors of the Federal Reserve System (U.S.) Sep 2019.

Canay, Ivan A., “A simple approach to quantile regression for panel data,” *Econometrics Journal*, October 2011, 14 (3), 368–386.

Carriero, Andrea, Todd E. Clark, and Marcellino Massimiliano, “Nowcasting Tail Risks to Economic Activity with Many Indicators,” Working Papers 20-13R2, Federal Reserve Bank of Cleveland May 2020.

Chen, Chung-Chi, Yu-Lieh Huang, and Fang Yang, “Semantics matter: An empirical study on economic policy uncertainty index,” *International Review of Economics & Finance*, January 2024, 89 (Part A), 1286–1302.

Correa, Ricardo, Keshav Garud, Juan M. Londono, and Nathan Mislang, “Sentiment in Central Banks’ Financial Stability Reports,” *International Finance Discussion Papers*, March 2017, (1203).

Covitz, Daniel, Nellie Liang, and Tobias Adrian, “Financial Stability Monitoring,” *Annual Review of Financial Economics*, December 2015, 7 (1), 357–395.

Detken, Carsten, Olaf Weeken, Lucia Alessi, Diana Bonfim, Miguel M. Boucinha, Christian Castro, Sebastian Frontczak, Gaston Giordana, Julia Giese, Nadya Jahn, Jan Kakes, Benjamin Klaus, Jan Hannes Lang, Natalia Puzanova, and Peter Welz, “Operationalising the countercyclical capital buffer: indicator selection, threshold identification and calibration options,” ESRB Occasional Paper Series No. 5, ESRB June 2014.

Engle, Robert F. and Susana Campos-Martins, “What are the events that shake our world? Measuring and hedging global COVOL,” *Journal of Financial Economics*, None 2023, 147 (1), 221–242.

Fell, John and Garry Schinasi, “Assessing Financial Stability: Exploring the Boundaries of Analysis,” *National Institute Economic Review*, 2005, 192, 102–117.

- , **Sandor Gardo, Benjamin Klaus, Jonas Wendelborn, and Stefan Wredenborg**, “Communication for financial crisis prevention: a tale of two decades,” *ECB Financial Stability Review*, November 2024, 2.
- Fernandez-Villaverde, Jesus, Tomohide Mineyama, and Dongho Song**, “Are We Fragmented Yet? Measuring Geopolitical Fragmentation and Its Causal Effect,” NBER Working Papers 32638, National Bureau of Economic Research, Inc Jun 2024.
- Figueres, Juan Manuel and Marek Jarocinski**, “Vulnerable growth in the euro area: Measuring the financial conditions,” *Economics Letters*, 2020, 191 (C).
- Gentzkow, Matthew, Bryan Kelly, and Matt Taddy**, “Text as Data,” *Journal of Economic Literature*, September 2019, 57 (3), 535–574.
- Giacomini, Raffaella and Ivana Komunjer**, “Evaluation and Combination of Conditional Quantile Forecasts,” *Journal of Business & Economic Statistics*, October 2005, 23, 416–431.
- Kremer, Manfred, Marco Lo Duca, and Daniel Hollo**, “CISS - a composite indicator of systemic stress in the financial system,” Working Paper Series 1426, European Central Bank Mar 2012.
- Kwon, Byeunchun, Taejin Park, Phurichai Rungcharoenkitkul, and Frank Smets**, “Parsing the Pulse: Decomposing Macroeconomic Sentiment with LLMs,” *BIS Working Papers*, October 2025, (1294).
- Lang, Jan Hannes, Cosimo Izzo, Stephan Fahr, and Josef Ruzicka**, “Anticipating the bust: a new cyclical systemic risk indicator to assess the likelihood and severity of financial crises,” Occasional Paper Series 219, European Central Bank February 2019.
- , **Marek Rusnák, and Moritz Greiwe**, “Medium-Term Growth-at-Risk in the Euro Area,” *IMF Economic Review*, 2025.
- Lefort, Baptiste, Eric Benhamou, Jean-Jacques Ohana, David Saltiel, Beatrice Guez, and Damien Challet**, “Can ChatGPT Compute Trustworthy Sentiment Scores from Bloomberg Market Wraps?,” *SSRN Online Library, Paris Dauphine University*, February 2024, pp. 1–32.
- Lopez-Lira, Alejandro and Yuehua Tang**, “Can ChatGPT Forecast Stock Price Movements? Return Predictability and Large Language Models,” *arXiv preprint*, April 2023.

- Loughran, Tim and Bill McDonald**, “When Is a Liability Not a Liability? Textual Analysis, Dictionaries, and 10-Ks,” *The Journal of Finance*, January 2011, 66 (1), 35–65.
- Ludvigson, Sydney C., Sai Ma, and Serena Ng**, “Uncertainty and Business Cycles: Exogenous Impulse or Endogenous Response?,” *American Economic Journal: Macroeconomics*, October 2021, 13 (4), 369–410.
- Nyman, Rickard, Sujit Kapadia, and David Tuckett**, “News and narratives in financial systems: Exploiting big data for systemic risk assessment,” *Journal of Economic Dynamics & Control*, June 2021, 127.
- Schularick, Moritz and Alan M. Taylor**, “Credit booms gone bust: Monetary policy, leverage cycles, and financial crises, 1870-2008,” *American Economic Review*, April 2012, 102 (2), 1029–61.
- Schüler, Yves S., Paul P. Hiebert, and Tuomas A. Peltonen**, “Financial cycles: Characterisation and real-time measurement,” *Journal of International Money and Finance*, 2020, 100 (C).
- Zhang, Zijian, Rong Fu, Yangfan He, Xinze Shen, Yanlong Wang, Xiaojing Du, Haochen You, Jiazhao Shi, and Simon Fong**, “FinSentLLM: Multi-LLM and Structured Semantic Signals for Enhanced Financial Sentiment Forecasting,” *Preprint in arXiv*, September 2025.

Annex

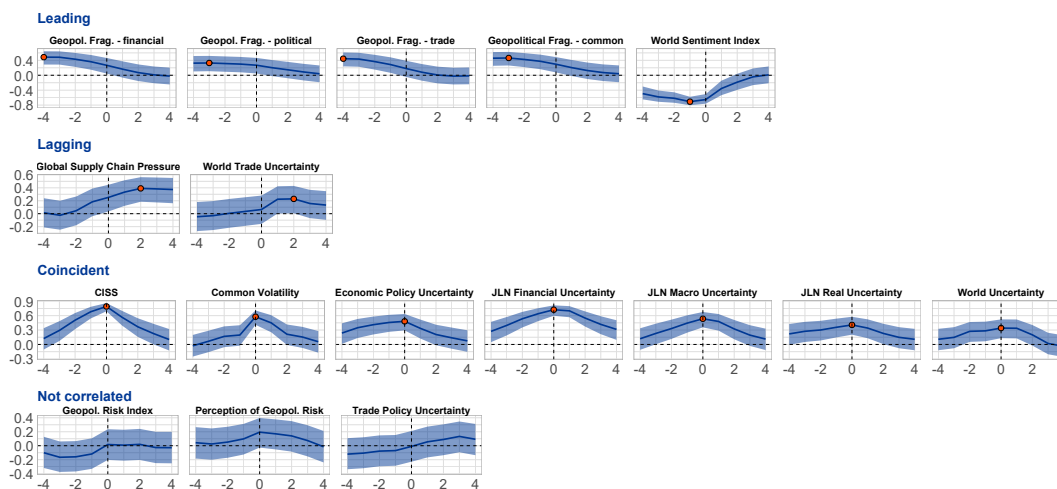
In the Annex we show additional tables and figures.

A.1 Additional data and indicator characteristics

Table A.1: Data and Classification Summary

| Dataset characteristics | Value |
|--|---------------------|
| Raw total number of articles | 1,068,423 |
| Sample period | Jan 2005 – Jan 2026 |
| Average article length (words) | 582 |
| Average article length (characters) | 4,151 |
| Duplicates removed | 13,661 |
| Classification pipeline | |
| Articles entering Stage 1 | 1,052,734 |
| Excluded in Stage 1 (non-economic content) | 408,972 (38.8%) |
| Articles entering Stage 2 | 643,762 |
| Excluded in Stage 2 (non-trigger events) | 499,585 (77.6%) |
| Articles entering Stage 3 | 143,897 |

Figure A.1: Lead-lag relationship between SPOT and other indicators



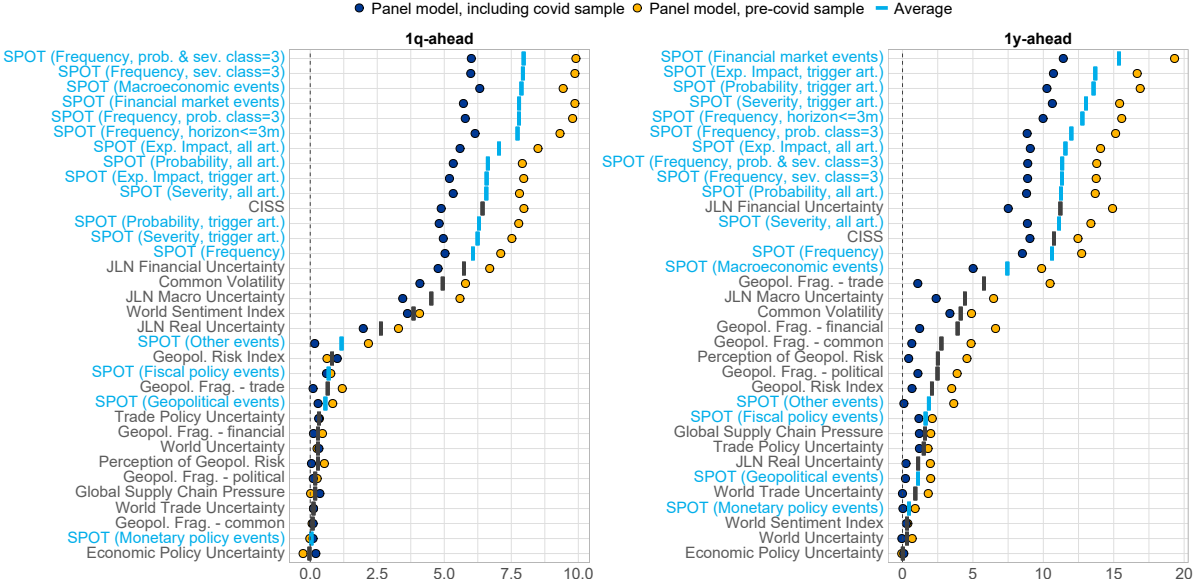
Notes: Computed on quarterly sample 2005Q1-2024Q4 using 3m averages of monthly values. Red dot denotes maximum correlation (significant at 5% level). Blue lines denote cross-correlations between SPOT - expected impact (all articles) and a given indicator at various leads (negative x-axis) and lags (positive x-axis). Light blue area denotes 95% confidence intervals.

A.2 Robustness of Growth-at-Risk results

Because various modelling choices might affect the key results above, we also estimate alternative specifications to assess their robustness. First, we test whether the results are sensitive

to the estimation sample. A key concern is that the large drops in real GDP during the COVID pandemic might disproportionately drive the results, even though the shock was exogenous. To address this, we include a set of dummy variables. As an additional check, we re-estimate the models on a sample that ends in 2019 Q4. [Figure A.2](#) reports these results. The yellow dots show the improvement from adding indicators relative to the baseline model estimated solely on the pre-COVID sample. The results indicate that, although including the COVID period reduces overall model performance somewhat, the relative improvement from the indicators is very similar to that in the model estimated on the full sample including the COVID period and COVID dummies.

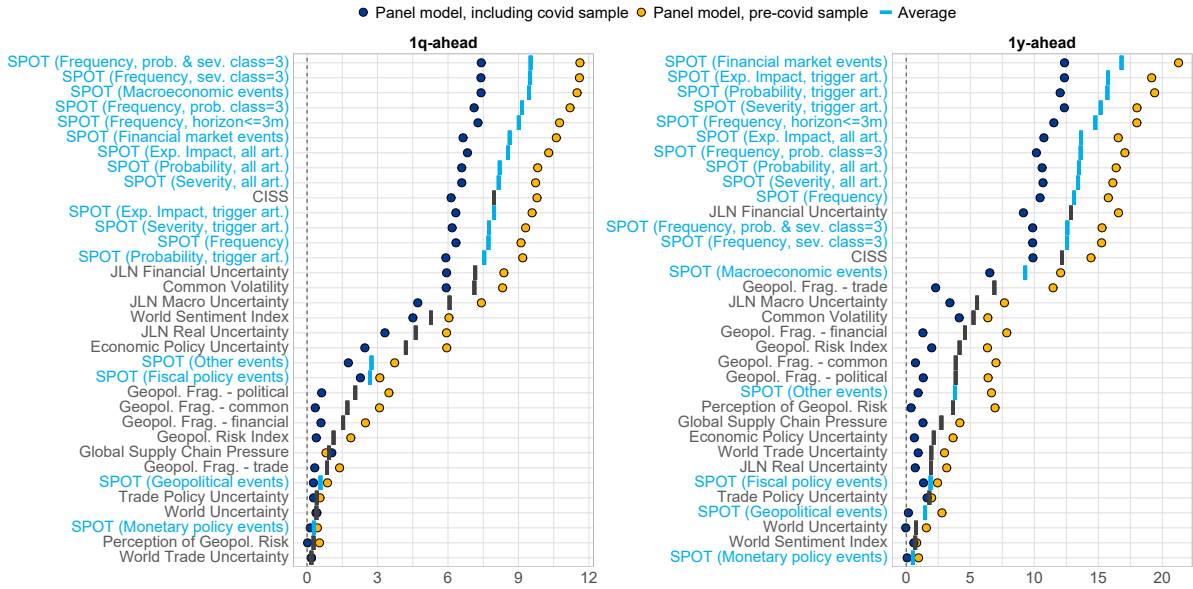
Figure A.2: Performance of the indicators for various sample periods



Notes: Improvements in tick loss function are in percentage relative to the baseline model with GDP and SRI. The models are ordered by the average performance of indicators on the two samples.

Similarly, [Figure A.3](#) compares models with interactions estimated on the sample that includes the COVID period (dark blue dots) with those estimated on the sample that ends before the pandemic (yellow dots). As with the models without interactions, the key findings of the benchmark specification continue to hold.

Figure A.3: Performance of the indicators with interactions for various sample periods



Notes: Improvements in tick loss function are in percentage relative to the baseline model with GDP and the SRI (including its lag and a dummy variable for cases when SRI>0). The models are ordered by the average performance of indicators on the two samples.

Next, we examine how sensitive the results are to the choice of quantile level. Setting $\tau = 0.1$ means that the estimates are driven by the lowest 10% of observations. While focusing on the lower tail is in line with the key goal of the study—namely, to capture rare, extreme crisis periods that are most relevant for financial stability surveillance—the number of influential observations is rather small. For example, for quarterly data over a 20-year period, the results are driven by only the eight lowest values for each country. We therefore re-estimate the models for the 25th quantile ($\tau = 0.25$), which still reflects the lower part of the GDP distribution but is based on more observations. These results, reported in Figure A.4, indicate that the key messages from the benchmark specification remain largely unchanged, despite the higher quantile, thereby corroborating the robustness of our findings.

Figure A.4: Performance for models evaluated at different quantile ($\tau = 0.25$)

