

Saufex

D3.2 – Report on AI-based Tools effectiveness

Table of contents

1. Introduction.....	3
2. Terms and Definitions.....	4
3. Recent Developments in Implemented AI Tools.....	6
3.1 Updates on LLM Usage.....	6
3.2 Customized AI Solutions.....	8
3.3 Future Prospects for Generative AI.....	10
4. Effectiveness Evaluation of Implemented AI Tools.....	13
4.1 TTP Classification.....	13
4.2 Custom TTP Classification.....	20
4.3 Narrative Detection.....	20
4.4 Custom Narrative Detection.....	21
4.5 Named Entity Recognition.....	22
4.6 Custom Named Entity Recognition.....	25
5. Planned Updates and Improvements.....	26
6. Conclusions.....	27
7. References.....	28

1. Introduction

This report builds on [Deliverable 1.3 – Report on Existing AI-based Tools](#) and presents updated findings on selected AI tools, an evaluation of their implementation effectiveness, and the next steps for improving the solution.

The information environment is evolving at a blistering pace. Foreign Information Manipulation and Interference (FIMI) campaigns now harness deep-fakes, cross-platform “narrative laundering,” and highly targeted micro-influence operations that mutate week by week. At the same time, AI capabilities – especially large language models (LLMs) and agent-based orchestration frameworks – have leapt forward, promising analysts richer insight, faster triage, and unprecedented linguistic reach. The pressing question is no longer whether AI can contribute to FIMI defence, but which tools work best for which tasks, how reliably, and at what cost.

Deliverable D1.3 reviewed the first wave of commercially and academically available AI systems. That assessment identified a clear pattern: small, fine-tuned classifiers performed best on well-defined, slowly changing tasks, whereas GPT-class LLMs were stronger in reasoning-heavy or fast-evolving domains – provided they were guided by well-designed prompts. It also anticipated a shift from single-step text processing toward agentic AI pipelines that can plan, retrieve information, call external APIs, and execute multi-step workflows to support analysts.

This report addresses three practical questions and describes the implemented application through that lens:

- How effectively do LLM-centred workflows perform on core FIMI-defence tasks – classifying tactics, techniques, and procedures (TTPs), detecting and normalising influence narratives, and extracting named entities from multilingual, noisy open-source text?
- Where do fine-tuned specialised models, retrieval components, or lightweight adapters add measurable value, and where is prompt engineering sufficient?
- How can we give subject-matter experts direct control – enabling them to upload custom TTP subsets, narrative lists, or entity glossaries – without requiring coding or model retraining?

2. Terms and Definitions

Foreign Information Manipulations and Interference (FIMI) refers to the actions taken by foreign entities to deliberately manipulate and interfere with the information environment of a target country or organization. These actions are intended to influence public opinion, disrupt societal cohesion, and undermine the integrity of democratic processes. FIMI encompasses a range of activities, including the dissemination of disinformation, misinformation, and propaganda.

Natural Language Processing (NLP) is a field of artificial intelligence that focuses on the interaction between computers and humans through natural language. NLP involves the application of computational techniques to analyze, understand, and generate human language, enabling machines to process and respond to text and speech data.

Named Entity Recognition (NER) is a Natural Language Processing (NLP) subtask that automatically identifies and classifies key information – such as names of people, organizations, locations, and dates – within unstructured text. It acts as an information extraction tool, transforming raw text into organized data for improved search, analysis, and content recommendation.

Natural Language Processing (NLP) is a field of artificial intelligence that focuses on the interaction between computers and humans through natural language. NLP involves the application of computational techniques to analyze, understand, and generate human language, enabling machines to process and respond to text and speech data.

Bidirectional Encoder Representations from Transformers (BERT) is a state-of-the-art NLP model developed by Google. BERT uses a transformer-based architecture that allows it to understand the context of words in a sentence bidirectionally, meaning it considers the entire sentence when interpreting each word. This capability makes BERT highly effective for a wide range of language understanding tasks.

Large Language Models (LLMs) are advanced NLP models that are trained on extensive datasets comprising vast amounts of text. Examples include GPT (Generative Pre-trained Transformer), Llama, Gemini, and Mistral. These models can generate human-like text and are capable of understanding and responding to complex language inputs.

Tactics, Techniques, and Procedures (TTPs) refer to the methods and strategies used by adversaries to achieve their objectives. In the context of FIMI, TTPs encompass the various approaches used to spread disinformation, manipulate narratives, and interfere with the information environment of a target.

Narrative is a structured story or account of connected events, whether real or imagined, typically featuring characters, a setting, and a conflict. It represents a specific way of organizing, explaining, and understanding experiences, acting as a core element of communication across literature, media, and daily life.

Generative AI describes algorithms (such as ChatGPT) that take unstructured data as input (for example, natural language and images) to create new content, including audio, code, images, text, simulations, and videos. It can automate, augment, and accelerate work by tapping into unstructured mixed-modality data sets to generate new content in various forms.

Applied AI technologies and techniques use models trained through machine learning to solve classification, prediction, and control problems in order to automate activities, add or augment capabilities and offerings, and improve decision making.

Agentic AI is an artificial intelligence system capable of independently planning and executing complex, multistep tasks. Built on foundation models, these agents can autonomously perform actions, communicate with one another, and adapt to new information. Significant advancements have emerged, from general agent platforms to specialized agents designed for deep research.

3. Recent Developments in Implemented AI Tools

In Deliverable D1.3 (Report on Existing AI-based Tools), drawing on sources [2–7], we concluded that model suitability depends on the task. Narrow, well-defined problems often benefit from specialised, fine-tuned ML models, while GPT-4-class LLMs tend to perform better on reasoning-intensive or less structured tasks. On fixed propaganda or misinformation benchmarks, RoBERTa-style baselines can outperform out-of-the-box LLMs unless prompts are carefully designed. However, adding task context, clear definitions, and few-shot examples can significantly reduce this gap.

For “open” or fast-changing tasks – such as emerging TTPs, new narratives, or evolving named entities – the flexibility of LLMs often makes them the more practical choice. In computational social science, LLMs already show strong performance across sentiment analysis, hate-speech detection, stance detection, misinformation classification, and event extraction, frequently in zero- or few-shot settings. Key limitations that remain relevant include bias, limited explainability, inference cost, sensitivity to prompting, and the ongoing arms race with evolving misinformation tactics.

3.1 Updates on LLM Usage

In the second half of 2025, new research papers [8–25] further confirmed the value of LLM-based approaches for the target tasks addressed in this deliverable.

In DiNaM: Disinformation Narrative Mining with Large Language Models [8], Sosnowski et al. (2025) introduce a four-stage pipeline that automates the discovery of disinformation narratives in fact-checking corpora. Its core contribution is the structured use of GPT-class LLMs. Across the evaluated models (GPT-4o-mini, Qwen3-32B, Llama-3-70B, Gemma-3-27B), GPT-based systems consistently achieved the strongest results. The authors also highlight that strong zero- and few-shot performance reduces the need for costly task-specific fine-tuning – aligning with our own implementation experience. Structured prompting (role definition, explicit guidelines, and a strict output schema) enables GPT to operate as an end-to-end information extraction and synthesis component rather than a simple summariser. When evaluated using a geometry-aware metric (WCD), GPT-based modules outperform specialised non-generative systems in both claim extraction and narrative synthesis. Overall, the findings suggest that carefully engineered prompts can turn general-purpose LLMs into reliable and cost-efficient tools for large-scale disinformation analysis – without new training data or bespoke model development [8].

More broadly, the 2025 literature indicates increasing adoption of LLMs in FIMI-relevant toolchains. Key drivers include:

- **Broad linguistic and domain coverage “out of the box”**

Open-weight or API-based LLMs have already been shown to match or exceed specialist models on eight benchmark datasets drawn from politics, health and science without any fine-tuning [9]. Additionally, the researches underline three main factors contributing to worse performance of LLMs: (a) failure to effectively integrate analysis results, (b) wrong analytical process, and (c) unrelated analyses.

Retrieval-augmented LLMs such as RAEoLLM boost cross-domain accuracy by up to 31 percentage points compared with conventional few-shot baselines, while requiring no task-specific training [10]. Because most frontier LLMs are pre-trained on dozens of languages, they can detect or summarise manipulation across linguistic boundaries that graph-based or monolingual classifiers miss; PolyTruth and X-Troll report the best F1 scores on low-resource Slavic and African-language troll data when a multilingual LLM backbone is used [11].

- **Faster adaptation through zero-/few-shot prompting rather than costly re-training**

Analysts across institutions (including EEAS and NATO) note that FIMI narratives can shift on a weekly cadence [12]. Prompt-based workflows allow new indicators and definitions to be introduced immediately, whereas graph- or SVM-based pipelines typically require additional feature engineering and iteration. In An Agentic Operationalization of DISARM for FIMI Investigation on Social Media [14], Tseng et al. (2026) describe a DISARM-aligned multi-agent pipeline in which specialised LLM agents map previously unseen manipulation behaviours to the DISARM taxonomy in near real time, scaling tasks that historically required manual triage.

- **Built-in explainability options that increase analyst trust**

Token-level explanations (e.g., SHAP/LIME-style rationales) derived from LLM outputs correlate with human fact-checker rationales, providing an auditable trail for operational use and reducing the risk that automated flags are dismissed [15].

- **Operational scalability and cost profile**

Once deployed, an LLM with retrieval or lightweight adapters can analyse millions of posts per hour on commodity GPUs, whereas graph-propagation models require laborious crawl-and-merge steps. DISARM-agent experiments processed two weeks of Twitter data (≈ 120 M posts) in < 7 hours on a 4xA100 node [14]. Cloud-hosted APIs (e.g., GPT-4o-mini) now price at a fraction of a cent per 1k tokens, under-cutting the engineering and labeling costs of maintaining separate models for each threat theatre. In the D1.3 – Report on Existing AI-based Tools, we already mentioned prediction of LLMs costs lowering during development of new models and competition growth.

Late-2025/early-2026 studies increasingly show that LLMs deliver higher accuracy, faster adaptation, richer explanations and better multi-modal coverage than earlier keyword, SVM or

GNN-only tools in the FIMI domain. When coupled with retrieval, adapter or agentic scaffolding, LLMs not only detect hostile narratives but also help analysts understand, simulate and proactively neutralise them – making LLMs the strategic core of next-generation cognitive-security stacks.

Specifically in TTP classification, there are new solutions using LLMs. For example, TRIAGE combines two LLM modules (rule-prompted + in-context) to map fresh CVEs to techniques. The hybrid LLM raises recall over rule-only mapping and shows that GPT-4o-mini already outperformed Llama-3.3-70B five months after both models were released. Feedback-driven Instruction Refinement (TTPrompt) lets the same prompts evolve automatically when MITRE renames or splits techniques, a task that used to take CTI teams days of manual rule editing [16].

In narrative detection, recent papers argue that LLMs capture narrative structure more deeply, including story grammar, causal chains, and frame semantics. This enables recognition of narratives even when relevant keywords do not co-occur – for instance, linking text about “renewable subsidies killing jobs” to a broader “climate-policy criticism” narrative [18]. A framing-theory model that injects frame elements into an LLM detects “re-framed” misinformation far better than feature-engineered baselines, confirming the value of narrative – semantic knowledge [19]. Other advantages mentioned in latest research papers are higher accuracy – often without any task-specific training, cross-lingual generalisation, faster adaptation to emergent or evolving narratives [20-23]. At the same time, LLMs may hallucinate or mislabel niche, newly emerging narratives – one of the reasons we introduced custom narrative detection capabilities.

Named Entity Recognition remains one of the most widely studied tasks, while still posing substantial practical challenges. Recent work reports similar LLM advantages as for TTP and narrative detection: strong zero- and few-shot performance approaching (and sometimes surpassing) supervised baselines, multilingual and fine-grained entity coverage, improved generalisation to unseen or emerging entities, and better handling of complex structures such as nested/overlapping entities and long documents [24]. Our own implementation experience also suggests that some established solutions (e.g., GLiNER) can be less flexible and less accurate than LLM-based approaches while requiring higher computational resources.

3.2 Customized AI Solutions

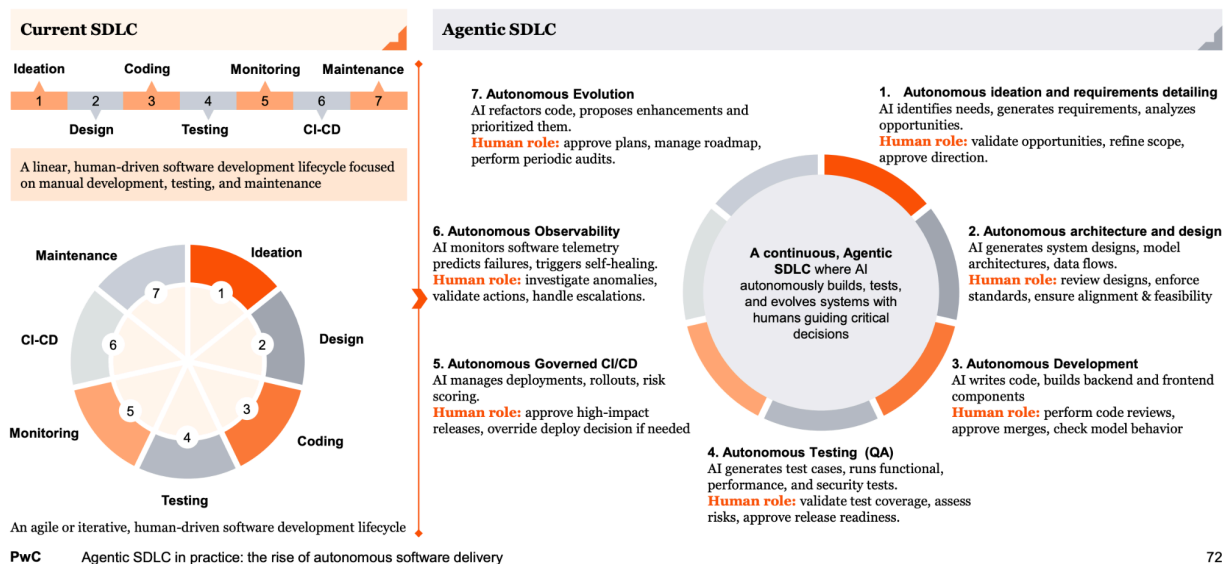
One of the most notable AI trends is the emergence of agentic use cases built on GPT-class models. According to McKinsey’s report “The State of AI in 2025: Agents, Innovation, and Transformation” [26], 62 % of respondents are already testing “AI agents” – autonomous GPT-based workflows that can plan, call outside tools, and trigger downstream actions. 23 % say they are scaling at least one agentic system; however, penetration by function is low (≤ 10 % in any single domain). Early traction is in IT service desks, knowledge-management research, marketing content generation, and software-engineering copilots. These agents are typically

built by wrapping GPT-4-class models with orchestration layers (LangChain, Microsoft AutoGen, CrewAI) plus function-calling APIs.

McKinsey itself developed QuantumBlack, an AI practice inside McKinsey & Company. Originating in Formula 1 where real-time telemetry demanded rapid model-driven decisions, it now functions as the firm's "AI consulting arm," mixing sector strategy teams with data-science and ML engineering squads. Its declared edge is "hybrid intelligence": combining human domain expertise with large-scale generative models.

Ernst&Young also embeds "leading-edge AI capabilities" throughout its services by reusing data and components from the firm's EY Fabric platform [31]. BCG helps organizations combine predictive AI and generative AI, "weaving together human and technological capabilities" to achieve large-scale productivity, cost, and innovation gains [32].

Completing the "Big Four" set, PwC's report *Agentic SDLC in Practice: The Rise of Autonomous Software Delivery* [27] focuses on Agentic SDLC – A software-delivery lifecycle where autonomous or oversighted AI agents plan, code, test, deploy and operate features with minimal human intervention, guided by high-level intent. They introduced the future Agentic SDLC concept, which they believe is a likely direction for future software development.



72

Beyond enterprise deployment, Hasselwander et al. (2026), in *Toward Agentic AI: User Acceptance of a Deeply Personalized AI Super Assistant (AISA)* [28], suggest that consumer readiness for deeply-personalized agent AI has moved from cautious interest to confident expectation – provided systems stay enjoyable, useful, easy, and trustworthy. Meeting those four conditions while managing perceived risk will shape whether AISAs (AI Super-Assistant) become the dominant interface to digital life in the immediate future.

Overall, current research and industry practice indicate a move toward more interoperable and specification-driven AI systems: agents that are less dependent on a single model version, easier to govern through explicit schemas and constraints, more interactive with users, and more capable of handling new task types and object classes as operational requirements evolve.

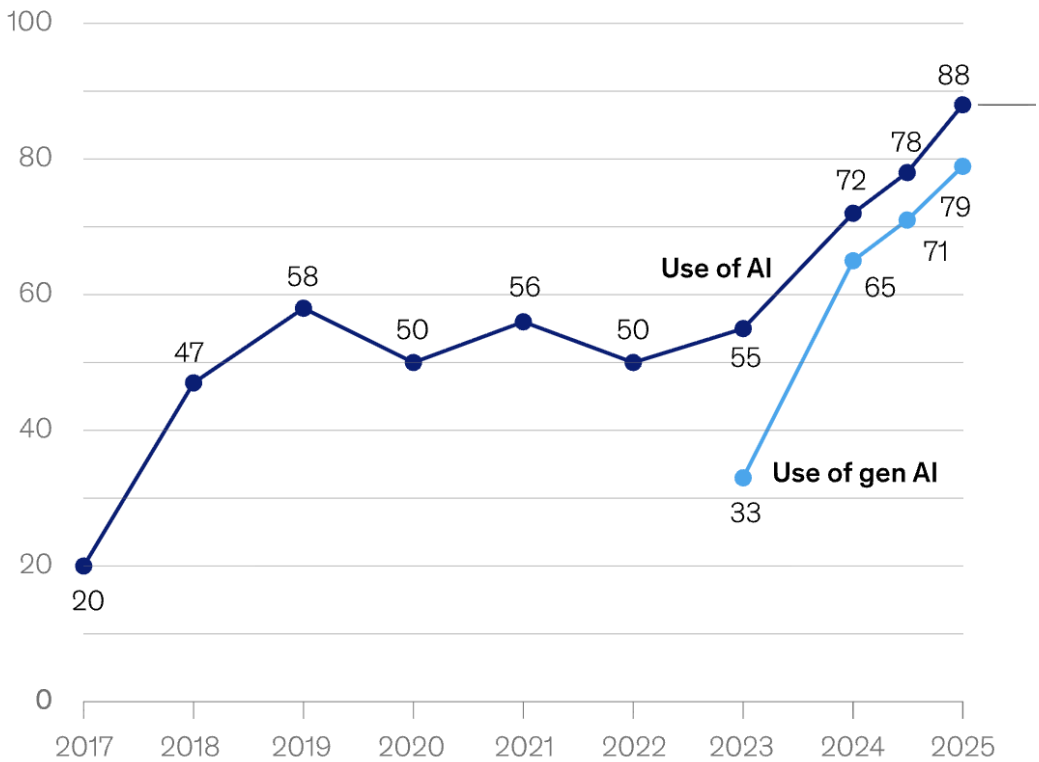
3.3 Future Prospects for Generative AI

According to McKinsey's The State of AI in 2025: Agents, Innovation, and Transformation (November 2025) [26], AI has become mainstream: 88 % of the 1,993 executives surveyed (105 countries, all sectors) report that their companies now use AI in at least one business function, up from 78 % in 2024. Generative-AI models – predominantly GPT-style large language models – are the engine of this growth. Most deployments sit on top of commercial APIs such as OpenAI GPT-4o, Anthropic Claude, Google Gemini, or on an enterprise-tuned variant of these models (Azure OpenAI, AWS Bedrock, etc.).

The impact of AI and GenAI on headcount remains ambiguous. Over the next 12 months, 32 % of respondents expect net workforce reductions ≥ 3 %, 13 % expect increases, 43 % no change. Larger enterprises foresee bigger reductions, but simultaneously continue to hire AI specialists.

In Deliverable D1.3, we used McKinsey's AI adoption indices to illustrate that AI usage – especially generative AI – was accelerating across industries. Updated numbers for 2025 suggest an even faster adoption trajectory than in 2024.

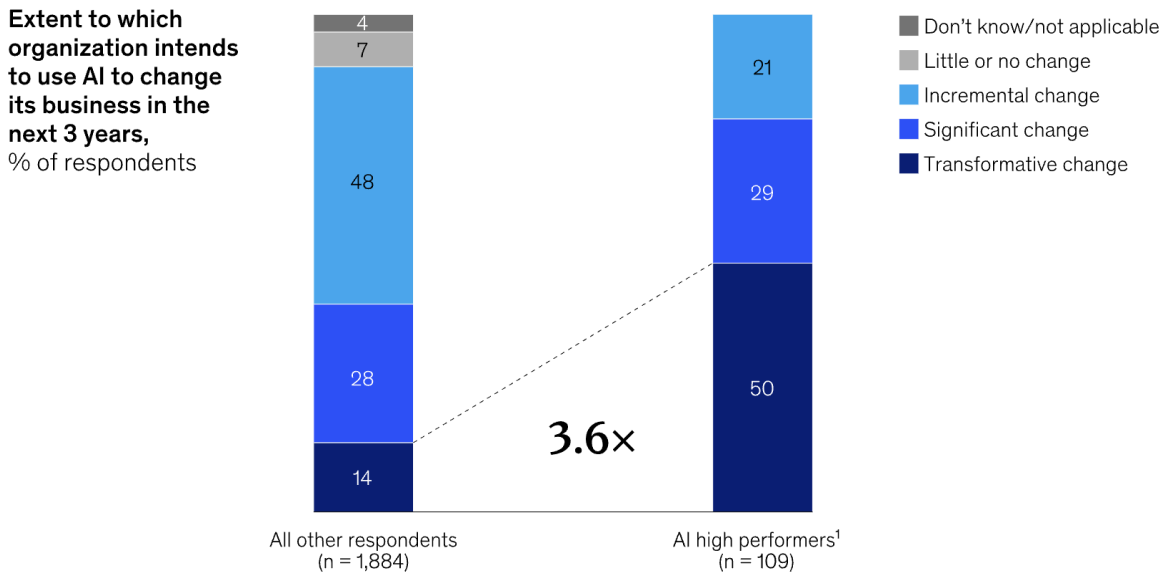
Use of AI by respondents' organizations, % of respondents
Organizations that use AI in at least 1 business function¹



Source: McKinsey. "The State of AI in 2025: Agents, Innovation, and Transformation". November 2025

It is also indicative that high performers express more ambitions to transform their business: AI high performers are more than three times more likely than others to say their organization intends to use AI to bring about transformative change to their businesses.

High performers are more likely than others to expect their organizations to use AI for enterprise-wide transformative change.



Source: McKinsey. “The State of AI in 2025: Agents, Innovation, and Transformation”. November 2025

Also, high performers are also nearly three times as likely as others are to say their organizations have fundamentally redesigned individual workflows.

At the same time, risk mitigation patterns remain uneven: inaccuracy is one of the two risks most commonly addressed by organisations, whereas explainability – the second-most frequently reported risk – is not among the most commonly mitigated [26].

External industry reporting reinforces this picture of rapid diffusion. Forbes characterises GenAI adoption as “one of the fastest in tech history,” noting it has outpaced cloud adoption by roughly 40% [29].

PwC similarly finds that optimism grows with experience: more mature teams are 10–20 percentage points more likely to expect deeper integration, improved training, and higher levels of automation from GenAI [27]. PwC’s survey-based projection suggests that by 2027 more than half of Middle East software teams will operate a fully augmented SDLC, rising to roughly two-thirds by 2029 – a pattern that may be indicative of wider global trends in software delivery automation [27].

Year	Assumed cumulative Pioneer share	Rationale
2025 (survey)	38 %	Baseline measurement
2026	46 %	8-point uplift driven by 76 % “Likely” investors
2027	54 %	Early-majority threshold crossed
2029 (median)	65 %	Syncs with Deloitte/Gartner agentic milestone
2030	~70 %	Plateau as late adopters close gap

Source: *Agentic SDLC in practice: the rise of autonomous software delivery*. PwC, 2026

Finally, Woollacott (2026), reporting on Gartner’s outlook on “AI model collapse” and data reliability, notes that by 2028 the proliferation of unverified AI-generated data could prompt 50% of organisations to adopt a zero-trust posture for data governance. This trend highlights an important implication: even as AI systems become more capable, the demand for strong human expertise – particularly in validation, governance, and analytical judgement – may increase rather than diminish.

By contrast, **EU official stats for 2025** show a more conservative but *rapidly increasing* footprint of AI in business operations, with substantial variation by firm size and sector – giving a useful **regional balance and evidence base** for EU-focused reports. The EU has moved from early experimentation to the early-majority phase of enterprise AI, propelled primarily by generative-language applications. Growth is fast but not yet fast enough to meet 2030 ambitions, and it remains highly uneven across firm size, sectors and Member States. Closing the gap will depend less on new algorithms than on human-capital development, SME-friendly infrastructure, and clear, innovation-compatible regulation [33] Reuters notes, to hit 2030 ambitions, the need is to pair generative-AI enthusiasm with capability-building, regulatory clarity and core-digital uptake [34].

4. Effectiveness Evaluation of Implemented AI Tools

To complement the three core tasks – TTP classification, narrative detection, and named entity recognition (NER) – we implemented customisation options for each module. This was motivated by the rapidly evolving operating environment: new narratives emerge daily, taxonomies and TTP definitions are updated over time, and the set of relevant entities expands continuously. Custom options give users the flexibility to focus on specific objects and adapt the system to new monitoring needs as the information space becomes more dynamic and complex.


The underlying model is configurable and can be changed by the developer as needed. To reduce dependence on any single model family, we designed the prompts in a unified format and avoided model-specific assumptions wherever possible.

Below presented the current test application interface, which Osavul and its partners use to evaluate the effectiveness of the implemented functionality.

TTP Classification Custom TTP Classification Narrative Detection Custom Narrative Detection NER Detection Custom NER Detection

TTP Classification

Upload a PDF or DOCX file ?

 Drag and drop file here
Limit 200MB per file • PDF, DOCX

Browse files

As noted above, current AI systems are trending toward greater independence, stronger specification-driven behaviour, deeper user interaction, and the ability to handle new task types and object classes. These trends are directly reflected in operational user needs and were a key reason for implementing the additional custom functionality across all three tasks: TTP, narrative, and NER detection.

4.1 TTP Classification

The TTP Classification module identifies relevant tactics, techniques, and procedures (TTPs) in the input text. For each detected TTP, the module outputs the mapped TTP entry, the supporting text fragment, and a relevance score on a 1–9 scale.

This output is useful in two common scenarios:

1. When a complete list of TTPs present in a report is required, the full set of results can be reviewed and exported.
2. When a single best-matching TTP is needed for a specific fragment, users can sort by score and select the top-ranked item.

The output table includes the following fields: TTP identifier (full number), full TTP description, supporting text fragment, and relevance score. In the current implementation, results are returned for scores 6–9, where 9 indicates the highest correspondence and 6 indicates moderate correspondence.

After uploading a file, users receive a results table that can be resized, filtered, and downloaded.

TTP Classification

Custom TTP ClassificationNarrative DetectionCustom Narrative DetectionNER DetectionCustom NER Detection

TTP Classification

Upload a PDF or DOCX file

Drag and drop file here

Limit 200MB per file • PDF, DOCX

Browse files

2024 Incident Alert Report template_DTL.docx

34.7KB

Analysis Results

Full TTP number	Full TTP description	fragment
T0019	Flood social channels; drive traffic/engagement to all assets; create aura/sense/perce	The inauthentic account seeded 105 and amplified 90 posts. The p
T0023	Change, twist, or exaggerate existing facts to construct a narrative that differs from re	Meta Narrative: the US is evil. Sub-narratives: Smearing US presiden
T0084	When an operation recycles content from its own previous operations or plagiarizes f	The account plagiarized several political cartoons from the actual a
T0085	Creating and editing false or misleading text-based artifacts, often aligned with one o	x.com/FdsRene impersonated a Dutch cartoonist to amplify anti-U
T0085.003	An influence operation may develop false or misleading news articles aligned to their	The inauthentic account seeded 105 and amplified 90 posts. The p
T0086	Creating and editing false or misleading visual artifacts, often aligned with one or mo	The account plagiarized several political cartoons from the actual a
T0086.001	Memes are one of the most important single artefact types in all of computational pro	The posts mostly contained anti-US, anti-Western, and pro-China c

Report on AI-based Tools effectiveness

Full TTP number	Full TTP description	fragment	score
T0019	Flood social channels; drive traffic/engagement to all assets; create aura/sense/perception	The inauthentic account seeded 105 and amplified 90 posts. The posts mostly contained anti-US, anti-Western, and pro-China content. Some of this content was reposted or shared from Chinese state media outlets or government officials. The account was part of a network of at least 8 accounts that actively amplified content seeded by x.com/husband_s1. We believe this network can almost certainly be linked to Spamouflage.	7
T0023	Change, twist, or exaggerate existing facts to construct a narrative that differs from reality		8
T0084	When an operation recycles content from its own previous operations or plagiarizes		9
T0085	Creating and editing false or misleading text-based artifacts, often aligned with one or more specific narratives	x.com/FdsRene impersonated a Dutch cartoonist to amplify anti-US, anti-Western, and	8
T0085.003	An influence operation may develop false or misleading news articles aligned to their own narrative	The inauthentic account seeded 105 and amplified 90 posts. The posts mostly contain	7
T0086	Creating and editing false or misleading visual artifacts, often aligned with one or more specific narratives	The account plagiarized several political cartoons from the actual artist.	7
T0086.001	Memes are one of the most important single artefact types in all of computational propaganda	The posts mostly contained anti-US, anti-Western, and pro-China content. Some of the	6
T0007	Create key social engineering assets needed to amplify content, manipulate algorithms	The account was part of a network of at least 8 accounts that actively amplified content	8
T0090	Inauthentic accounts include bot accounts, cyborg accounts, sockpuppet accounts, and	Between 2023-01-10 and 2024-09-14, an inauthentic account (x.com/FdsRene) impersonated	9
T0090.004	Sockpuppet accounts refer to falsified accounts that either promote the influence operation's	The account was part of a network of at least 8 accounts that actively amplified content	8
T0098.002	Leverage Existing Inauthentic News Sites	The inauthentic account seeded 105 and amplified 90 posts. The posts mostly contain	7
T0099	An influence operation may prepare assets impersonating legitimate entities to further	Between 2023-01-10 and 2024-09-14, an inauthentic account (x.com/FdsRene) impersonated	9
T0102.001	Use existing Echo Chambers/Filter Bubbles	The account was part of a network of at least 8 accounts that actively amplified content	7
T0104	Social media are interactive digital channels that facilitate the creation and sharing of	x.com/FdsRene impersonated a Dutch cartoonist to amplify anti-US, anti-Western, and	8
T0104.001	Examples include Facebook, Twitter, LinkedIn, etc.	Both x.com/FdsRene and x.com/husband_s1 were suspended after Voice of America (7
T0045	Use the fake experts that were set up during Establish Legitimacy. Pseudo-experts are	x.com/FdsRene impersonated a Dutch cartoonist to amplify anti-US, anti-Western, and	7
T0116.001	Use government-paid social media commenters, astroturfers, chat bots (programmed	The account was part of a network of at least 8 accounts that actively amplified content	7
T0118	An influence operation may amplify existing narratives that align with its narratives to	The inauthentic account seeded 105 and amplified 90 posts. The posts mostly contain	9
T0128.002	Concealing network identity aims to hide the existence an influence operation's network	The account was part of a network of at least 8 accounts that actively amplified content	7
T0083	An influence operation may seek to exploit the preexisting weaknesses, fears, and enemies	Meta Narrative: the US is evil. Sub-narratives: Smearing US presidential candidates, p	8
T0115	Delivering content by posting via owned media (assets that the operator controls).	The inauthentic account seeded 105 and amplified 90 posts.	9
T0049	Flooding and/or mobbing social media channels feeds and/or hashtag with excessive	The account was part of a network of at least 8 accounts that actively amplified content	7

During early experiments, relying only on hyperlinks to the TTP framework or on the model's internal knowledge produced inconsistent mappings and, in some cases, hallucinated TTP references. To mitigate this, we added a structured TTP table as an explicit input to the analysis, improving grounding and reducing hallucinations.

TTPs considered of high and average importance and frequency in the tested reports

Full TTP number	Full TTP description
T0083	An influence operation may seek to exploit the preexisting weaknesses, fears, and enemies of the target audience for integration into the operation's narratives and overall strategy. Integrating existing vulnerabilities into the operational approach conserves resources by exploiting already weak areas of the target information environment instead of forcing the operation to create new vulnerabilities in the environment.
T0087	Creating and editing false or misleading video artifacts, often aligned with one or more specific narratives, for use in a disinformation campaign. This may include staging videos of purportedly real situations, repurposing existing video artifacts, or using AI-generated video creation and editing technologies (including deepfakes).
T0089.002	Create inauthentic documents intended to appear as if they are authentic non-public documents. These documents can be "leaked" during later stages in the operation.
T0115	Delivering content by posting via owned media (assets that the operator controls).

Report on AI-based Tools effectiveness

Full TTP number	Full TTP description
T0049	Flooding and/or mobbing social media channels feeds and/or hashtag with excessive volume of content to control/shape online conversations and/or drown out opposing points of view. Bots and/or patriotic trolls are effective tools to achieve this effect.
T0119	Cross-posting refers to posting the same message to multiple internet discussions, social media platforms or accounts, or news groups at one time. An influence operation may post content online in multiple communities and platforms to increase the chances of content exposure to the target audience.
T0066	Plan to degrade an adversary's image or ability to act. This could include preparation and use of harmful information about the adversary's actions or reputation.
T0003	Use or adapt existing narrative themes, where narratives are the baseline stories of a target audience.
T0004	Advance competing narratives connected to the same issue ie: on one hand deny incident while at same time expresses dismiss.
T0022	"Conspiracy narratives" appeal to the human desire for explanatory order, by invoking the participation of powerful (often sinister) actors in pursuit of their own political goals. These narratives are especially appealing when an audience is low-information, marginalized or otherwise inclined to reject the prevailing explanation. Conspiracy narratives are an important component of the "firehose of falsehoods" model.
T0022.001	An influence operation may amplify an existing conspiracy theory narrative that aligns with its incident or campaign goals. By amplifying existing conspiracy theory narratives, operators can leverage the power of the existing communities that support and propagate those theories without needing to expend resources creating new narratives or building momentum and buy in around new narratives.
T0019	Flood social channels; drive traffic/engagement to all assets; create aura/sense/perception of pervasiveness/consensus (for or against or both simultaneously) of an issue or topic.
T0023	Change, twist, or exaggerate existing facts to construct a narrative that differs from reality. Examples: images and ideas can be distorted by being placed in improper content.
T0084	When an operation recycles content from its own previous operations or plagiarizes from external operations. An operation may launder information to conserve resources that would have otherwise been utilized to develop new content.
T0085	Creating and editing false or misleading text-based artifacts, often aligned with one or more specific narratives, for use in a disinformation campaign.
T0085.003	An influence operation may develop false or misleading news articles aligned to their campaign goals or narratives.
T0086	Creating and editing false or misleading visual artifacts, often aligned with one or more specific narratives, for use in a disinformation campaign. This may include photographing staged real-life situations, repurposing existing digital images, or using image creation and editing technologies.
T0086.001	Memes are one of the most important single artefact types in all of computational propaganda. Memes in this framework denote the narrow image-based definition. But that naming is no accident, as these items have most of the important properties of Dawkins' original conception as a self-replicating unit of culture. Memes pull together reference and commentary; image and narrative; emotion and message. Memes are a powerful tool and the heart of modern influence campaigns.

Report on AI-based Tools effectiveness

Full TTP number	Full TTP description
T0087.001	Deepfakes refer to AI-generated falsified photos, videos, or soundbites. An influence operation may use deepfakes to depict an inauthentic situation by synthetically recreating an individual's face, body, voice, and physical gestures.
T0087.002	Cheap fakes utilize less sophisticated measures of altering an image, video, or audio for example, slowing, speeding, or cutting footage to create a false context surrounding an image or event.
T0007	Create key social engineering assets needed to amplify content, manipulate algorithms, fool public and/or specific incident/campaign targets. Computational propaganda depends substantially on false perceptions of credibility and acceptance. By creating fake users and groups with a variety of interests and commitments, attackers can ensure that their messages both come from trusted sources and appear more widely adopted than they actually are.
T0013	Create media assets to support inauthentic organizations (e.g. think tank), people (e.g. experts) and/or serve as sites to distribute malware/launch phishing operations.
T0090	Inauthentic accounts include bot accounts, cyborg accounts, sockpuppet accounts, and anonymous accounts.
T0090.001	Anonymous accounts or anonymous users refer to users that access network resources without providing a username or password. An influence operation may use anonymous accounts to spread content without direct attribution to the operation.
T0090.004	Sockpuppet accounts refer to falsified accounts that either promote the influence operation's own material or attack critics of the material online. Individuals who control sockpuppet accounts also manage at least one other user account. Sockpuppet accounts help legitimize operation narratives by providing an appearance of external support for the material and discrediting opponents of the operation.
T0011	Hack or take over legitimate accounts to distribute misinformation or damaging content.
T0098.002	Leverage Existing Inauthentic News Sites
T0099	An influence operation may prepare assets impersonating legitimate entities to further conceal its network identity and add a layer of legitimacy to its operation content. Users will more likely believe and less likely fact-check news from recognizable sources rather than unknown sites. Legitimate entities may include authentic news outlets, public figures, organizations, or state entities.
T0018	Create or fund advertisements targeted at specific populations
T0101	Localized content refers to content that appeals to a specific community of individuals, often in defined geographic areas. An operation may create localized content using local language and dialects to resonate with its target audience and blend in with other local news and social media. Localized content may help an operation increase legitimacy, avoid detection, and complicate external attribution.
T0102.001	Use existing Echo Chambers/Filter Bubbles
T0104	Social media are interactive digital channels that facilitate the creation and sharing of information, ideas, interests, and other forms of expression through virtual communities and networks.
T0104.001	Examples include Facebook, Twitter, LinkedIn, etc.
T0105.002	Examples include Youtube, TikTok, ShareChat, Rumble, etc

Full TTP number	Full TTP description
T0045	Use the fake experts that were set up during Establish Legitimacy. Pseudo-experts are disposable assets that often appear once and then disappear. Give "credibility" to misinformation. Take advantage of credential bias
T0116.001	Use government-paid social media commenters, astroturfers, chat bots (programmed to reply to specific key words/hashtags) influence online conversations, product reviews, web-site comment forums.
T0118	An influence operation may amplify existing narratives that align with its narratives to support operation objectives.
T0119.001	An influence operation may post content across groups to spread narratives and content to new communities within the target audiences or to new target audiences.
T0119.002	An influence operation may post content across platforms to spread narratives and content to new communities within the target audiences or to new target audiences. Posting across platforms can also remove opposition and context, helping the narrative spread with less opposition on the cross-posted platform.
T0128.002	Concealing network identity aims to hide the existence an influence operation's network completely. Unlike concealing sponsorship, concealing network identity denies the existence of any sort of organization.

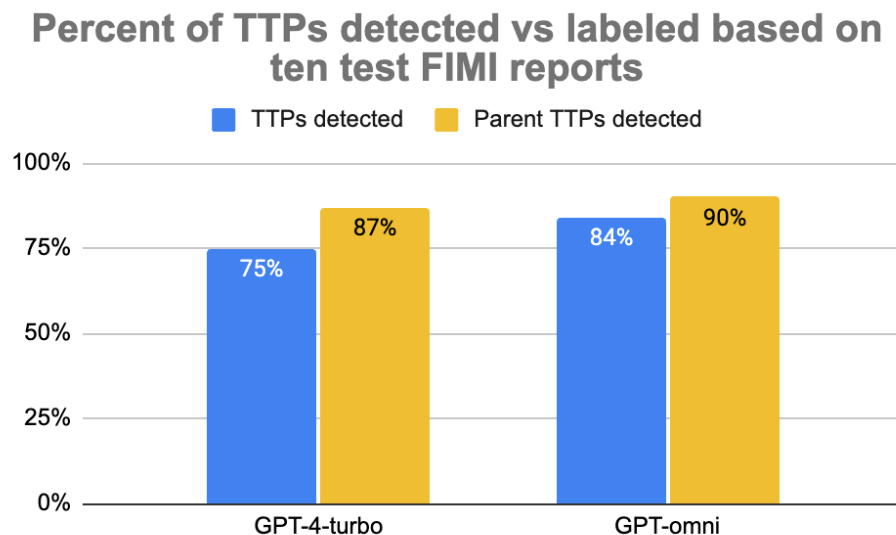
After a series of experiments, we selected the following approach as the best-performing configuration at this stage:

1. Match all the reports to all manually found TTPs, forming pairs of report text - TTP description.
2. For each pair, use LLM prompts to define:
 - report text fragment which corresponds to the TTP description
 - score of the correspondence from 0 to 10: 0 means that the found fragment does not correspond to the TTP, 5 - partly corresponds, 10 - completely corresponds
3. Define a score threshold above which the pair is considered a valid match and the TTP is treated as present in the report.
4. Evaluate performance across different parameter settings (e.g., thresholds, chunk sizes, overlap, and prompt variants) and compare results.

During Step 2, we initially tested a simplified prompt that only attempted to extract a supporting fragment and returned an empty output if no evidence was found. In practice, this produced mostly empty results. We therefore updated the prompt to always return a score alongside the fragment, enabling threshold-based selection and more stable behaviour.

For very long reports that exceed the effective context window or degrade model performance (partly due to overly broad input), we apply a preprocessing step: splitting reports into smaller overlapping chunks to preserve continuity and avoid cutting relevant evidence across boundaries.

As a result, we obtained an initial estimate of TTP-classification performance. For example, on a test set of ten FIMI reports, GPT-4o (“omni”) produced better results than earlier GPT-4-class models and, in our usage, offered a materially lower per-token cost (approximately 3.5-4× lower).



The underlying model remains configurable and can be swapped by the developer as LLM performance and pricing evolve quickly.

We also provided our SAUFEX project partner, Debunk, with access to the application and will incorporate their feedback to improve this module and the broader toolset. The most important planned improvements include:

- adding fields such as TTP name, TTP phase, and justification/rationale for analyst convenience;
- tightening the definition of TTP–evidence correspondence to reduce borderline matches;
- improving the output format to support a one-fragment–to-one-TTP display option (where appropriate).

4.2 Custom TTP Classification

Custom TTP Classification allows users to provide a custom list of TTPs with descriptions as an input parameter. The module then restricts detection to only those TTPs from the supplied list.

This capability is important for three reasons:

1. **Focused analysis:** analysts may need to monitor a specific subset of TTPs relevant to a particular case, threat actor, or campaign.
2. **Taxonomy evolution:** TTP definitions and taxonomies change over time, so user-managed lists can be updated without modifying the underlying system logic.
3. **Modular subsets:** users can select non-overlapping or weakly related TTP subsets to reduce ambiguity and improve interpretability of the results.

TTP Classification **Custom TTP Classification** Narrative Detection Custom Narrative Detection NER Detection Custom NER Detection

Custom TTP Classification

Upload a PDF or DOCX file



Drag and drop file here

Limit 200MB per file • PDF, DOCX

Browse files

Upload TTP list (PDF or DOCX file)



Drag and drop file here

Limit 200MB per file • PDF, DOCX

Browse files

4.3 Narrative Detection

The module extracts narratives from the input text and returns the detected narratives together with the supporting text fragments. In the current implementation, a narrative is defined as:

“A narrative is a deliberately constructed story that shapes perceived reality by selecting, emphasising, and framing information in a way that influences how people think, feel, and ultimately act.”


In practice, the exact definition of a “narrative” may vary by use case. For example, analysts may need to focus specifically on destructive narratives or narratives tied to a particular topic or threat context. For this reason, we treat the narrative definition as a configurable parameter where necessary; at the same time, the Custom Narrative Detection option


addresses a large part of this need by allowing users to constrain detection to a predefined narrative set.

Based on partner testing, we also introduced an explicit constraint to enforce unit narratives (“one narrative – one thought”). Concretely, narratives are required to be expressed as a single sentence containing exactly one core proposition. If multiple propositions are present, the output must be split into separate narratives. This change reduced the number of compound or mixed outputs and improved consistency by producing a list of atomic, reusable sub-narratives aligned to distinct text fragments.


TTP Classification Custom TTP Classification **Narrative Detection** Custom Narrative Detection NER Detection Custom NER Detection

Narrative Detection

Upload a PDF or DOCX file 

 Drag and drop file here
Limit 200MB per file • PDF, DOCX

Browse files

 Saufex Month 5.1 Annex 7 (Yaroslav).docx 0.7MB

×

Analysis Results

fragment_text	narrative_text
We found messages circulating claiming that during a phone call, Trump issued an ul	Trump issued an ultimatum to Zelensky, condemning him for corrupt ties to the Bide
This campaign continues to develop the narrative of the total corruption of the Ukrai	The campaign develops the narrative of the total corruption of the Ukrainian governr

Based on Debunk’s testing feedback, the following improvements are planned:

- Normalisation: move from raw claims or headline-like statements toward a normalised, reusable sub-narrative form suitable for cross-case comparison and trend tracking.
- Deduplication and consolidation: unify narratives that express the same underlying idea (e.g., multiple variations portraying Ukraine’s leadership as weak, fearful, or illegitimate) into a single stable sub-narrative, with variants mapped to that canonical form.

4.4 Custom Narrative Detection

Custom Narrative Detection enables users to search for specific, predefined narratives by providing a narrative list as an input parameter. The module then scans the input text and returns matches only from the supplied list, along with supporting text fragments.

Custom Narrative Detection

Upload a PDF or DOCX file

Drag and drop file here

Limit 200MB per file • PDF, DOCX

Browse files

Upload Narrative list (PDF or DOCX file)

Drag and drop file here

Limit 200MB per file • PDF, DOCX

Browse files

In testing with a custom list of 200+ narratives, the module demonstrated consistent retrieval across the full list – detecting narratives located at the beginning, middle, and end of the input list (i.e., without a bias toward early entries). In the example below, the first detected narratives originate from different sections of the same long custom list.

Custom Narrative Detection

Upload a PDF or DOCX file

Drag and drop file here

Limit 200MB per file • PDF, DOCX

Browse files

upd_Saufex Month 5.1 Annex 7 (Yaroslav).docx

2.8MB

×

Upload Narrative list (PDF or DOCX file)

Drag and drop file here

Limit 200MB per file • PDF, DOCX

Browse files

custom_narratives_example.docx

22.8KB

×

Analysis Results

fragment_text	narrative_text
This campaign continues developing the narrative of the total corruption of the Ukrai	S007.12 - Ukraine is a corrupt state / The Ukrainian government is corrupt.
As well as Ukraine doesn't follow international agreements supposed to be.	S007.4 - Ukraine fails to follow international agreements.
It was added that COVID-19 was definitely created as a bioweapon, which was unrela	S011.5 - COVID-19 was being engineered as a bioweapon.


4.5 Named Entity Recognition


Named Entity Recognition (NER) is one of the most widely used NLP tasks, yet it remains challenging in practice despite its seemingly straightforward definition. In our context, NER

covers three linked steps: extracting entity mentions from raw text, disambiguating them, and – where applicable – linking them to structured knowledge representations.



To address these requirements, we implemented a multi-step prompt workflow while keeping key components configurable. The most important inputs are a curated entity list with definitions and few-shot examples, which improve consistency across domains and languages. We also introduced a canonical entity name to merge different surface forms that refer to the same entity (e.g., “USA” and “United States”), improving normalisation and downstream analysis.

NER Detection

Upload REPORT (PDF or DOCX file) 

 Drag and drop file here
Limit 200MB per file • PDF, DOCX

Browse files

 2024 Incident Alert Report template_DTL.docx 34.7KB 

Analysis Results

Start	End	Entity	General_Entity	Label
1	1	2023-01-10	2023-01-10	Date
3	3	2024-09-14	2024-09-14	Date
9	9	x.com/FdsRene	x.com/FdsRene	Website
12	14	Dutch political Cartoonist	Dutch political Cartoonist	Identity
15	17	Bart van Leeuwen	Bart van Leeuwen	Persone
19	19	X	X	Media Channel
43	43	x.com/husband_s1	x.com/husband_s1	Website
55	55	Spamouflage	Spamouflage	Threat Actor
58	58	x.com/FdsRene	x.com/FdsRene	Website
60	60	x.com/husband_s1	x.com/husband_s1	Website

↓

🔍

⌵

Start	End	Entity	General_Entity	Label
1	1	2023-01-10	Date	Date
3	3	2024-09-14	Date	Date
10	10	x.com/FdsRene	x.com/FdsRene	Website
12	14	Dutch political Cartoonist	Bart van Leeuwen	Identity
16	16	Bart van Leeuwen	Bart van Leeuwen	Persone
18	18	X	X	Media Channel
47	47	x.com/husband_s1	x.com/husband_s1	Website
51	51	Spamouflage	Spamouflage	Threat Actor
55	55	x.com/FdsRene	x.com/FdsRene	Website
57	57	x.com/husband_s1	x.com/husband_s1	Website
60	62	Voice of America (VOA)	VOA	Media Channel
68	68	2024-09-19	Date	Date
72	72	x.com/FdsRene	x.com/FdsRene	Website
78	78	Global Times	Global Times	Media Channel
80	81	Chinese government official	Zhang Heqing	Identity
83	83	Zhang Heqing	Zhang Heqing	Persone
95	95	Spamouflage	Spamouflage	Threat Actor
97	97	DTL	DTL	Identity
99	99	VOA	VOA	Media Channel
101	101	2024-09-12	Date	Date
103	103	VOA	VOA	Media Channel

4.6 Custom Named Entity Recognition

Custom Named Entity Recognition extends the standard NER module by allowing users to provide a custom list of entities (with definitions/descriptions) as an input parameter. The module then scans the input text and returns matches restricted to the supplied list, improving precision when analysts need to monitor a specific set of entities.

An example of usage is shown below. In practice, this functionality can also be applied beyond classic NER: by shaping the input descriptions appropriately, users can search for specific text patterns or concept-defined mentions, not only named entities. This supports a broader and potentially more future-proof approach – enabling flexible expert-defined inputs that guide AI modules without requiring code changes or model retraining.

TTP ClassificationCustom TTP ClassificationNarrative DetectionCustom Narrative DetectionNER DetectionCustom NER Detection

Custom NER Detection

Upload REPORT (PDF or DOCX file)

Drag and drop file here
Limit 200MB per file • PDF, DOCX

Browse files

2024 Incident Alert Report template_DTL.docx34.7KB

X

Upload NER list (PDF or DOCX file)

Drag and drop file here
Limit 200MB per file • PDF, DOCX

Browse files

custom_ners_example.docx14.0KB

X

Analysis Results

Start	End	Entity	General_Entity	Label
68	68	x.com/husband_s1	x.com/husband_s1	Website
77	77	Spamouflage	Spamouflage	Threat Actor
81	81	x.com/FdsRene	x.com/FdsRene	Website
83	83	x.com/husband_s1	x.com/husband_s1	Website
87	88	Voice of America	VOA	Media Channel
90	90	VOA	VOA	Media Channel

5. Planned Updates and Improvements

Based on the research and the analysis of the implemented AI tools described above, our next iteration will focus on strengthening the core functionality; the custom modules will benefit from these improvements accordingly.

For TTP modules the following updates will be prioritized:

- Increase flexibility of TTP classification (e.g., better handling of overlapping/related techniques and edge cases).
- Improve result presentation and usability, including a more convenient output view.
- Add additional output fields for analyst use: TTP name, TTP phase, and justification/rationale.
- Tighten the definition of TTP–evidence correspondence to reduce ambiguous matches.
- Introduce an optional one-fragment–to-one-TTP display mode (where appropriate).

For narrative modules we will focus on:

- Move from raw claims or headline-like outputs toward a normalised, reusable sub-narrative format.
- Consolidate semantically equivalent narratives into stable canonical sub-narratives (e.g., multiple variations portraying Ukraine’s leadership as weak, fearful, or illegitimate).

We will also work on adding new customized parameters (e.g. narrative definition, as operationalised in DiNaM by Sosnowski et al., 2025 [8], NERs few-shot examples, etc.).

Technical updates based on our partner’s feedback will also include:

- Performance optimisation, especially improved processing speed for larger files.
- Support for additional input formats and a direct text input option (not only file upload).
- Translation options to support multilingual workflows.

6. Conclusions

The current Deliverable shifts from cataloguing available AI to using, optimising and operationalising it.

The AI tools development overview show several important tendencies which were considered during AI modules development and will influence its updates.

The performance-cost curve for foundation models has shifted dramatically. New model releases deliver higher accuracy at one-quarter to one-tenth of last year's inference price, while retrieval-augmented and few-shot prompting techniques routinely close the gap with specialist, fine-tuned models. As a result, LLMs have become the default choice for fast-moving, reasoning-heavy FIMI tasks – including TTP mapping, narrative discovery and multilingual NER – where their zero-/few-shot adaptability now outstrips classical SVM, GNN or RoBERTa baselines.

The ecosystem is moving from single-call text generators to agentic AI pipelines. Early “agent mode” deployments (ChatGPT, AutoGen, CrewAI) and consultancy frameworks (McKinsey's QuantumBlack, PwC's Agentic SDLC) show that autonomous chains of GPT calls can plan, retrieve, execute and refine tasks with minimal human intervention. These architectures already power IT service desks, CTI enrichment, content operations and software-engineering copilots, and they align with the need for high-throughput, continuously updated FIMI monitoring.

Flexibility for domain experts is now built in. Every core module – TTP, narrative, NER – can ingest custom taxonomies or entity lists on the fly, letting analysts steer the models without retraining. This “expert-in-the-loop” design mitigates hallucination risk, speeds response to emerging threats, and future-proofs the platform as taxonomies evolve.

Looking ahead, market signals confirm the trend. McKinsey finds 88 % of firms already use AI in at least one function; PwC predicts a majority of software teams will adopt fully agentic SDLC workflows by 2027; Gartner warns that unverified synthetic data will force half of all organisations into zero-trust governance by 2028. Taken together, these data points imply that the shift toward LLM-centric customisable tooling is not merely timely – it is essential for keeping pace with both the threat landscape and enterprise AI adoption curves. At the same time, the regional context is important. EU enterprise adoption of AI jumped to 20 % in 2025 – driven by generative-language tools and concentrated in large firms and tech-intensive sectors – but must accelerate roughly two-fold, especially among SMEs and lagging Member States, to meet the EU's 75 % Digital-Decade target by 2030.

The built functionality is analysed from the perspective of AI development, as well as user engagement and flexibility options. In partnership with Debunk, using testing results, the scope of improvements is proposed for the prioritized implementation.

7. References

1. Natural Language Inference Benchmarks and papers (<https://paperswithcode.com/task/natural-language-inference>)
2. Jose J., Greenstadt R. Are Large Language Models Good at Detecting Propaganda? New York University, Department of Computer Science and Engineering (2025). <https://arxiv.org/pdf/2505.13706>
3. Tianyi Huang, Jingyuan Yi, Peiyang Yu, Xiaochuan Xu. Unmasking Digital Falsehoods: A Comparative Analysis of LLM-Based Misinformation Detection Strategies. <https://arxiv.org/pdf/2503.00724>
4. Sosnowski W., Modzelewski A., Skorupska K., Otterbacher J., Wierzbicki A. EU DisinfoTest: a Benchmark for Evaluating Language Models' Ability to Detect Disinformation Narratives. 2024. <https://aclanthology.org/2024.findings-emnlp.862.pdf>
5. Thapa, S., Shiwakoti, S., Shah, S.B. et al. Large language models (LLM) in computational social science: prospects, current state, and challenges. Soc. Netw. Anal. Min. 15, 4 (2025). <https://doi.org/10.1007/s13278-025-01428-9>
6. Mingqi Gao, Xinyu Hu, Xunjian Yin, Jie Ruan, Xiao Pu, Xiaojun Wan. LLM-based NLG Evaluation: Current Status and Challenges, May 2025. https://direct.mit.edu/coli/article/doi/10.1162/coli_a_00561/128807/LLM-based-NLG-Evaluation-Cu
[rrent-Status-and](https://direct.mit.edu/coli/article/doi/10.1162/coli_a_00561/128807/LLM-based-NLG-Evaluation-Cu)
7. Singla A., Sukharevsky A., Yee L., Chui M., Hall B., March 2025. The state of AI. How organizations are rewiring to capture value. McKinsey, March 2025. https://www.mckinsey.com/~media/mckinsey/business%20functions/quantumblack/our%20insights/the%20state%20of%20ai/2025/the-state-of-ai-how-organizations-are-rewiring-to-capture-value_fin
[al.pdf](https://www.mckinsey.com/~media/mckinsey/business%20functions/quantumblack/our%20insights/the%20state%20of%20ai/2025/the-state-of-ai-how-organizations-are-rewiring-to-capture-value_fin)
8. Witold Sosnowski, Arkadiusz Modzelewski, Kinga Skorupska, and Adam Wierzbicki. 2025. [DiNaM: Disinformation Narrative Mining with Large Language Models](https://arxiv.org/abs/2505.13706). In *Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing*, pages 30224–30251, Suzhou, China. Association for Computational Linguistics. <https://aclanthology.org/2025.emnlp-main.1537.pdf>
9. Explore the Potential of LLMs in Misinformation Detection: An Empirical Study. Mengyang Chen, Lingwei Wei, Han Cao, Wei Zhou, Songlin Hu. December 2024. <https://arxiv.org/pdf/2311.12699>
10. RAEmoLLM: Retrieval Augmented LLMs for Cross-Domain Misinformation Detection Using In-Context Learning Based on Emotional Information. Zhiwei Liu, Kailai Yang, Qianqian Xie, Christine de Kock, Sophia Ananiadou, Eduard Hovy. May 2025. <https://arxiv.org/pdf/2406.11093>
11. X-Troll: eXplainable Detection of State-Sponsored Information Operations Agents. Lin Tian, Xiuzhen Zhang, Maria Myung-Hee Kim, Jennifer Biggs, Marian-Andrei Rizoiu. August 2025. <https://arxiv.org/pdf/2508.16021>
12. EEAS 3rd Threat Report. March 2025. <https://www.eeas.europa.eu/sites/default/files/documents/2025/EEAS-3nd-ThreatReport-March-2025-05-Digital-HD.pdf>
13. NATO StratCom COE. Virtual Manipulation Brief 2024/1. https://stratcomcoe.org/publications/download/VM_210x2975_FINAL_DIGITAL_PDF.pdf
14. An Agentic Operationalization of DISARM for FIMI Investigation on Social Media. Kevin Tseng, Juan Carlos Toledano, Bart De Clerck, Yuliia Dukach, Phil Tinn. January 2026. <https://arxiv.org/pdf/2601.15109>
15. Explaining Misinformation Detection Using Large Language Models. Vishnu S. Pendyala, Christopher E. Hall. 2024. <https://www.mdpi.com/2079-9292/13/9/1673>
16. A Systematic Approach to Predict the Impact of Cybersecurity Vulnerabilities Using LLMs. Anders Mølmen Høst, Pierre Lison, Leon Moonen. October 2025. <https://arxiv.org/pdf/2508.18439>

17. From Retrieval to Reasoning: A Framework for Cyber Threat Intelligence NER with Explicit and Adaptive Instructions. Jiaren Peng, Hongda Sun, Xuan Tian, Cheng Huang, Zeqing Li, Rui Yan. December 2025. <https://arxiv.org/pdf/2512.19414>
18. GateNLP at SemEval-2025 Task 10: Hierarchical Three-Step Prompting for Multilingual Narrative Classification. Iknoor Singh, Carolina Scarton, Kalina Bontcheva. May 2025. <https://arxiv.org/pdf/2505.22867>
19. Detecting misinformation through Framing Theory: the Frame Element-based Model. Guan Wang, Rebecca Frederick, Jinglong Duan, William Wong, Verica Rupar, Weihua Li, Quan Bai. February 2024. <https://arxiv.org/pdf/2402.15525>
20. COGNAC at SemEval-2025 Task 10: Multi-level Narrative Classification with Summarization and Hierarchical Prompting. Azwad Anjum Islam & Mark A. Finlayson. August, 2025. <https://aclanthology.org/2025.semeval-1.190.pdf>
21. Irapuarani at SemEval-2025 Task 10: Evaluating Strategies Combining Small and Large Language Models for Multilingual Narrative Detection. Gabriel Assis, Livia de Azevedo, João Vitor de Moraes, Laura Alvarenga and Aline Paes. August 2025. <https://aclanthology.org/2025.semeval-1.7.pdf>
22. GateNLP at SemEval-2025 Task 10: Hierarchical Three-Step Prompting for Multilingual Narrative Classification. Iknoor Singh, Carolina Scarton, Kalina Bontcheva. May 2025. <https://arxiv.org/pdf/2505.22867>
23. Narrative Shift Detection: A Hybrid Approach of Dynamic Topic Models and Large Language Models. Kai-Robin Lange, Tobias Schmidt, Matthias Reccius, Henrik Müller, Michael Roos and Carsten Jentsch. June 2025. <https://arxiv.org/pdf/2506.20269>
24. GPT-NER: Named Entity Recognition via Large Language Models. Shuhe Wang, Xiaofei Sun, Xiaoya Li, Rongbin Ouyang, Fei Wu, Tianwei Zhang, Jiwei Li, Guoyin Wang. May 2025. <https://aclanthology.org/2025.findings-naacl.239.pdf>
25. Structure-Aware Decoding Mechanisms for Complex Entity Extraction with Large-Scale Language Models. Zhimin Qiu, Di Wu, Feng Liu, Chenrui Hu, Yuxiao Wang. December 2025. <https://arxiv.org/pdf/2512.13980>
26. The state of AI in 2025: Agents, innovation, and transformation. McKinsey, November, 2025. <https://www.mckinsey.com/capabilities/quantumblack/our-insights/the-state-of-ai>
27. Agentic SDLC in practice: the rise of autonomous software delivery. PwC, 2026. <https://www.pwc.com/m1/en/publications/2026/docs/future-of-solutions-dev-and-delivery-in-the-rise-of-gen-ai.pdf>
28. Toward agentic AI: User acceptance of a deeply personalized AI super assistant (AISA). Marc Hasselwander, Varsolo Sunio, Oliver Lah, Emmanuel Mogaji. 2026. <https://www.sciencedirect.com/science/article/pii/S0969698925003996>
29. Generative AI's Unprecedented Adoption Cycle. Forbes, July 2025. <https://www.forbes.com/sites/timbajarin/2025/07/22/generative-ais-unprecedented-adoption-cycle>
30. Fears over "AI model collapse" are fueling a shift to zero trust data governance strategies. <https://www.itpro.com/security/data-protection/fears-over-ai-model-collapse-are-fueling-a-shift-to-zero-trust-data-governance-strategies>
31. Ernst&Young: AI Technology: Services & Solutions. https://www.ey.com/en_in/services/ai/technology
32. BCG. Artificial Intelligence at Scale. <https://www.bcg.com/capabilities/artificial-intelligence>
33. Use of artificial intelligence in enterprises. Eurostat. https://ec.europa.eu/eurostat/statistics-explained/index.php?title=Use_of_artificial_intelligence_in_enterprises
34. European small businesses rush into AI without basic digital tools, study shows. Reuters. <https://www.reuters.com/business/european-small-businesses-rush-into-ai-without-basic-digital-tools-study-shows-2025-10-08>