

Technical Specification Double-side Page

- 1. TECHNICAL SCOPE:** Summarize the mock-up devised during the EXPLORE phase: how have you addressed the challenge/Theme Challenges and tackled with its requirements and data. Include a diagram.

The lack of data flow and deep learning models between press agencies and media outlets hence incurs data silos in the media business. Breaking these silos using Federated learning approach, and creating the corresponding data value chains represents a great opportunity for the media business across Europe.

(Requirement 1) Provide to press agencies a service that offers to them ML-based methods to evaluate news text quality and predict news text performance;

(R2) Ensure media outlets metrics data privacy and security;

(R3) Enable transparency about entity's contribution (trusted);

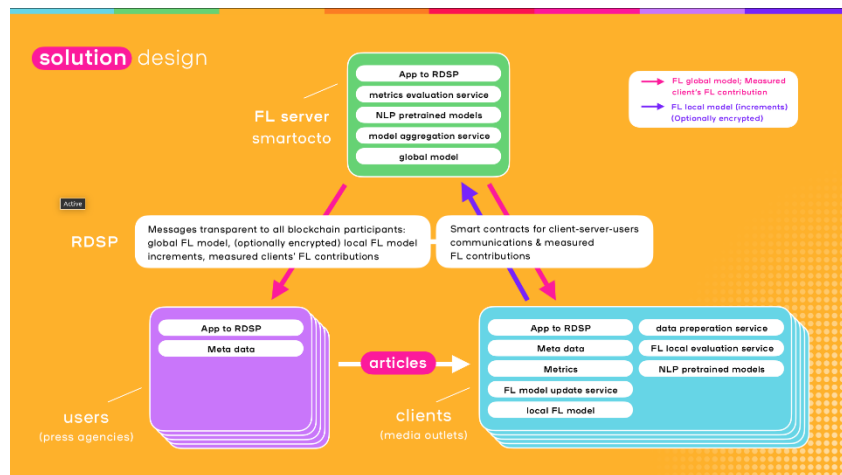
(R4) The service should be scalable with respect to data and model sizes;

(R5) The service should improve PA news text performance on a longer time scale;

We will use three datasets from two external partners, and one from the REACH data catalogue. The first dataset will be provided by the Dutch Press Agency. The second data set will be provided by the PA Media, a British press agency. In addition to these two external data sets, Smartocto will use a data set from the REACH catalogue – the VRT Dutch News Articles.

Smartocto acts as a FL service provider, while media outlets act as data providers (FL clients), and press agencies act here as users, getting access to globally trained FL models that they did not have before.

The blockchain and smart contract mechanisms, realized through the REACH data sharing platform (RDSP), ensure privacy and security of communications, as well as transparency and trustworthiness of the process.



- 2. ALGORITHMS, TOOLS AND CONCLUSIONS:** Detail the algorithms and tools identified to accomplish the challenge/Theme Challenges. Show clear understanding of the used REACH dataset/s and addressed challenge/Theme Challenges.

The FLAME approach will be compliant with state-of-the-art solutions for NLP like BERT, DistilBERT and BART that will serve as, pre-trained NLP models. FLAME utilizes FL and personalization to improve those models with a small computational and communication cost incurred. In terms of FL algorithms, FedAvg, FedProx, FedSGD, and FedOPT will be considered. In terms of open-source tools, besides RDSP, open-source tools for FL like FATE, LEAF, Tensorflowfederated, OpenFed, and for blockchain like Ethereum will be considered.



3. **SCALABILITY AND FLEXIBILITY OF THE SOLUTION:** Discuss whether the solution can truly cope with humongous and increasing datasets and how flexible it is to adapt to other related domains

The relevant aspects with respect to which scalability may be evaluated include: the trained ML model size; the data available; and the number of FL clients. With respect to model size, the FL system scalability can be ensured by utilizing algorithms that perform only partial model updates at each round. We will use pre-trained state-of-the-art model approaches such as BERT or even reduced-size models such as DistilBERT. Regarding data, scalability can be ensured by performing mini-batch type updates, where each local update corresponds to, e.g., a few rounds of (mini-batch) stochastic gradient descent. Regarding the number of clients, the FL algorithms with partial clients' participation resolve the issue. All the above FL mechanisms can be combined with the utilized blockchain and smart contracts-based approach and maintain scalability.

The envisioned FL supervised ML algorithms like FedAvg are generally applicable to multiple application domains, as well as the FL approach and the blockchain-based communications utilized.

4. **DATA GOVERNANCE AND LEGAL COMPLIANCE:** Describe the security level of the proposed solution, i.e. how authentication, authorization policies, encryption or other approaches are used to keep data secure. Explain how will be compliant with the current data legislations concerning security and privacy (e.g. GDPR).

The system uses the Reach data sharing platform from the Reach toolbox to facilitate blockchain and smart contracts-based operation. The server, clients (media outlets), and users (press agencies) have dApps that listen to the events on blockchain, hence realizing secure communication between actors. Several smart contracts are initiated on the blockchain. The first smart contract corresponds to implementing communication, and the second smart contract corresponds to recording each client's contribution in a transparent and immutable way. The communications that take place via blockchain include the following. The server receives from each client their current local models (or local model updates-increments). The clients and the users receive from the server the current global model. As communication on the blockchain is of broadcast-type, if sharing local models across different clients is an issue, then local models (or their increments) are encrypted via an asymmetric encryption protocol before being published on the blockchain, so that only the server sees the decrypted model updates. On the other hand, each client's contribution, at each FL iteration round, is published on the blockchain and visible to other actors.

In the development process we will use a privacy-by-design approach. Smartocto is compliant with EU data privacy legislation <https://smartocto.com/gdpr/>. All the personal data we collecting, are immediately, in the first step of processing, rendered anonymous in such a manner that the data subjects are no longer identifiable, so that system further processes only anonymized data. For the purpose of FLAME project, considering the specifics of the project and used Federated learning approach, clients and users don't need to share any sensitive data with the service provider (Smartocto).

5. **QUALITY ASSURANCE AND RISK MANAGEMENT:** Describe the quality process planed for the final product. Technologically, which are the potential risks in all the phases of the project (design of the solution, development, testing, deployment...) and indicate mitigation plans to still fulfil the challenge/Theme Challenges and data provider requirements.

Smartocto follows a best practice for the design of the solution architecture, and a strict and systematic approach for the testing and monitoring of data and features quality and the quality of respective models.

A potential risk is different structure of data across different data providers. In FL, this may translate into lack of labels or lack of features. This can be mitigated via: 1) using semi-supervised instead of supervised learning; 2) data and labels imputation. Another risk may be scalability of the blockchain. However, due to low volume and frequency of communication, because of using the pre-trained model approach, this is unlikely to be an issue. If a risk is identified, an alternative open source blockchain may be considered.

