# FEDGPT

## EXPLORE | EXPERIMENT | EXECUTE

## What is FedGPT?

FedGPT is a fully-containerized and vendor agnostic generative AI accelerator that enables federal agencies to build solutions with best-of-breed proprietary and open-source foundation models.

**UI**: Customizable interface with pre-built prompt libraries, monitoring dashboards, and builder mode

**APPS**: Pre-built gen AI apps like Chat, Marketing, Code Assist, Data Insights; ability to Build Your Own

**MODELS**: Prebuilt connectors for 12 foundation models along with standard benchmark metrics

**DATA**: Flexible & vendor agnostic data storage; pre-built RAG implementation with common vector DBs

**INFRASTRUCTURE**: On-prem & cloud-agnostic GPUs, containerized microservices

## Why FedGPT?

FedGPT offers a set of accelerators for rapid, safe and secure experimentation with LLMs. Choose from Azure content moderation services, NVIDIA NeMo guardrails, and custom-built safeguards to validate model inputs and outputs. FedGPT complies with DISA application security and development security STIG and CISA secure software development standards, built with security by design to expedite ATO process. FedGPT includes:

✓ Secure sandbox support to run Gen AI experiments
✓ Benchmarks to evaluate & select fit for purpose models
✓ Prebuilt applications to expedite prototyping and deployment
✓ Dev team to build custom use cases

## Modular Components

**Model garden** with 12 out-of-the-box models & performance benchmarks to select the right model for each use case.

**Prebuilt applications** including Code Assist and RAG-implementation to expedite use case build & deployment.

**Functional UI** including a prompt library with 30 templates and dashboards to monitor spend & consumption.

**Robust safeguards** to validate user inputs, curate model outputs, & ensure end-to-end security.

**Secure Architecture** suitable for accelerated deployment in the most restrictive environments.

# Generative AI Experimentation with FedGPT

| EXPLORE | | Test frontier proprietary and open- source models. Identify best fit model for federal use cases. |
| --- | --- | --- |
| EXPERIMENT | | Conceptualize and build generative AI use cases that address mission & business needs. |
| EXECUTE | | Deploy and scale use cases in production considering sustainable solutions and support. |

Select an experimentation approach based on the desired use case fidelity. Note: prices exclude technology costs.

| PROTOTYPE | PROOF OF CONCEPT | PILOT |
| --- | --- | --- |
| ~$150K-300K | ~$400K-$600K | ~$1.1M-$1.3M |
| Test concept w/ FedGPT UI | Build MVP with custom UI | Deploy and scale |
| 6 to 8 weeks | 10 to 12 weeks | 4+ months |
| 4-person dev pod | 7-person dev pod | 10-person dev pod |

## Deploying FedGPT

Our teams will work with you to identify the deployment method best suited to your agency's requirements. Choose from the following:

| | Models | Data Type | Infra Needs | IT Costs | Scale |
| --- | --- | --- | --- | --- | --- |
| **AFS Cloud** | CSP APIs or open-source | Public data only | AFS procured | Token-based usage & compute costs | Prototype or POC |
| **Agency Cloud** | | Public + Client Data | Agency procured | | Prototype to prod deployment |
| **Agency On-Prem** | Open-source | | | GPU server hosting | |

## Ready to Get Started?

For more information and to get started today, please contact our FedGPT leads, Shauna and Conor.

Shauna Revay
Shauna.revay@afs.com

Conor McSherry
Conor.mcsherry@afs.com