

АРХИТЕКТУРА ЭКОСИСТЕМЫ ГЕНЕРАТИВНОГО ИСКУССТВЕННОГО ИНТЕЛЛЕКТА

Версия v.1.0

Генеративные ИИ-приложения и сервисы

Горизонтальные и вертикальные прикладные решения, доступные одновременно по модели on-premise и в облачном развертывании

Горизонтальные приложения

Финансы BloombergGPT FinanceGPT J.P.Morgan	Обслуживание клиентов Gerwin radaar	Продажи Microsoft AI SEARCH Apolloio DataRobot	Продуктивность copy.ai RoboGPT Notion	Информационные технологии GitHub Copilot J.NNs AI CRUDERRA	Управление знаниями character.ai Notion Deept
---	---	---	--	--	--

Вертикальные/отраслевые приложения

Финтех НейроТекстер Copilot	Естественные науки Recursion NVIDIA BioNeMo Morphic Labs	Здравоохранение zakupoint Health IBM ARTERA	PR/маркетинг CopyMonkey YandexGPT 2 Jasper AI Tome Николай Иронов Adobe Firefly	Юриспруденция Humata LEGAISE DECODER alawyer	Ритейл Meta	Транспорт GHOST	Робототехника Google
-----------------------------------	---	--	--	---	----------------	--------------------	-------------------------

Инструменты для генеративного ИИ

Инструменты для эффективного и безопасного промышленного использования базовых моделей генеративного ИИ

Кейс-специфичные

- Платформы для чат-ботов: Replika, Cleverbots, InflectionPi, alta
- Чат-боты (сочетание классических подходов с LLM): RASA, botpress
- Агенты: Jay Copilot, imbue, Just AI
- Транскрибация, суммаризация и перевод текста: M.T.S AI, Vision, happyscribe, DeepScribe, YandexGPT 2, cohere
- Генеративные презентеры: synthesia, D-ID
- Генерация видео: GEN-2, Sora, LUMIERE

Модель-центричные

- Разработка модели: GitHub, OpenAI, Google colab, Jupyter
- Риск-ориентированное использование моделей: LexCheck, Timely, Carebrix
- Retrieval Augmented Generation (RAG): Microsoft, Azure AI Search, Just AI, Confident AI, Jay Copilot
- Файн-тьюнинг: Jay Copilot, Just AI, M.T.S AI

Дата-центричные

- Промпт инжиниринг: Spellbook, vellum, Promptly
- Векторные базы данных и эмбединги: Pinecone, drant, Weaviate, chroma ai, milvus

Модели генеративного ИИ

GPT-модели saiga_mistral, YandexGPT 2, Gemini, ANTHROPIC	Диффузионные модели Stable Diffusion, Kandinsky 3.0, Шердурм, stability.ai	Доменные модели AI PICASSO, ChefGPT, Николай Иронов, BloombergGPT	Маркетплейсы генеративных моделей cloud.ru, Yandex Cloud, Hugging Face
---	---	--	---

Подготовка и обмен данными

- Федеративное обучение: Google, IBM
- Обезличивание: DATPROF, ataccama, DEEP DOCS
- Гомоморфное шифрование: Google, Microsoft
- Инструментарий для формирования датасетов: sweetviz, Superset
- Синтетические данные: hazy, kview, D, Dotonize, scale, surge, Synt Data
- Разметка данных: TAGME, Toloka, Label Studio

AI TRISM-инструменты

- Мониторинг отравления и искажения промптов: WHYLABS, Prompt
- Контроль галлюцинаций: LAKERA, Lasso, WHYLABS
- Контроль деанонимизации данных: LAKERA, CALYPSOAI, vigilai
- Модерация результатов (анти-байес): LLM GUARD, ROBERT INTELLIGENCE, vigilai
- Обеспечение безопасности ИИ-приложений: Lasso, Prompt, LLMfuzzer, LAKERA, SLASH-NEXT, vigilai, BurpGPT, LLM GUARD
- Объяснимость результатов ИИ: vigilai, LAKERA, BurpGPT
- Предотвращение кибератак с ИИ: Jericho Security, CALYPSOAI, BRIGHTSIDE, ZIRCH, ADVERSA, LLM GUARD

Модели генеративного ИИ

Модели и инструменты генеративного ИИ, реализующие базовые функции технологии

Технологические инструменты

Инструменты MLOps DataSphere, mlflow, Kubeflow, SinaraML, SELDON	Фронтенд для нейросетей stable-diffusion-webui, text-generation-webui	NoCode LowCode платформы Directual, appflowy, NOCODB, NocoBase	Оптимизация нейромоделей run:ai, deci, Skoltech	Оркестрация моделей по API Hubble, stack, Patterns
---	--	---	--	---

Инфраструктура

Техническая инфраструктура для обучения и применения (инференса) инструментов на основе генеративного ИИ

Инфраструктура

Вычисления для обучения NVIDIA, AMD, HUAWEI, intel, AMD, HUAWEI, intel, Google	Вычисления для интерфейса AMD, intel, HUAWEI	Программируемые/специализированные интегральные схемы AMD XILINX, LATTICE, FUJIAN MICROELECTRONICS GROUP, GOWIN	Сети связи Juniper, CISCO, NVIDIA, H3C, NETWORKS	Интерконнект NVIDIA, ИНТЕРКОННЕКТ, Ростелеком	Облачные решения CloudMTS, Selectel, Azure, AWS, Google Cloud	Анклавные вычисления (изолированные) BTB	Приватные вычисления (Privacy Enhanced Computation) Azure AI Search, Microsoft, СБЕР	Средства обеспечения базовой безопасности positive technologies, BIZZONE, КОА безопасности, kaspersky, TCC
---	---	--	---	--	--	---	---	---

Авторы

Максим Григорьев, Алексей Сидорук, Евгений Морозов, Валерия Митина, Варвара Ольшницкая

Контрибьютеры

Александр Долбнев (Яндекс), Полина Гришина (Сбер Девайсы), Феликс Скворцов (МТС), Василий Кирсанов (CloudAmplifier), Николай Бушков (Альфа Банк), Сергей Колесников (Тинькофф), Владислав Тушканов (Лаборатория Касперского), Ильяс Киреев (Positive Technologies)

