

**ПОЛТАВСЬКИЙ ДЕРЖАВНИЙ АГРАРНИЙ УНІВЕРСИТЕТ
НАВЧАЛЬНО-НАУКОВИЙ ІНСТИТУТ ЕКОНОМІКИ, УПРАВЛІННЯ,
ПРАВА ТА ІНФОРМАЦІЙНИХ ТЕХНОЛОГІЙ
КАФЕДРА ІНФОРМАЦІЙНИХ СИСТЕМ ТА ТЕХНОЛОГІЙ**

Пояснювальна записка

до кваліфікаційної роботи на здобуття ступеня вищої освіти Бакалавр
на тему: «Персональний асистент фермера на основі великих мовних
моделей»

Виконав: здобувач вищої освіти
за освітньо-професійною програмою
Інформаційні управляючі системи
спеціальності 126 Інформаційні
системи та технології
ступеня вищої освіти Бакалавр
групи 126ІСТбд41
Шубка М.І.
Керівник: Слюсар В.І.
Рецензент: Брикун О.М.

Полтава – 2024 року

ВСТУП

Актуальність теми кваліфікаційної роботи підтверджується необхідністю підвищити продуктивність, стійкість і прибутковість сільськогосподарських операцій. Для допомоги фермерам у прийнятті поінформованих рішень, плануванні діяльності та адаптації до мінливих умов та ін. можливо використовувати штучний інтелект. Він дозволить обробляти великі обсяги агрономічних даних, надавати навчальні матеріали, поради щодо лікування рослин, ідентифікувати хвороби або шкідники та ін. В сучасних умовах найбільшого поширення набув штучний інтелект на основі нейронних мереж. При цьому фермер може використовувати спеціальний додаток-асистент на мобільному пристрої, в тому числі, без необхідності підключення до Інтернету. В свою чергу, зараз спостерігається значний розвиток мереж на основі архітектури трансформер. Однак питання автоматизації роботи фермера за рахунок впровадження штучного інтелекту на основі великих мовних моделей (LLM) потребує додаткових досліджень. Все це свідчить про актуальність теми роботи.

Метою кваліфікаційної роботи є зменшення когнітивного навантаження на фермерів в процесах їхньої діяльності за рахунок інтелектуальних функцій на основі великих мовних моделей.

Завданнями кваліфікаційної роботи є:

- обґрунтування вибору інструментарію для створення асистента фермера;
- оцінка властивостей LLM;
- формування рекомендацій щодо реалізації асистента фермера на основі LLM;
- економічне обґрунтування прийнятих рішень.

Об'єктом дослідження є процес взаємодії фермерів з мовними моделями в контексті аграрної діяльності.

Предметом дослідження є особливості адаптації великих мовних моделей при рішенні специфічних завдань аграрного сектору.

Методами визначення інструментарію для розробки асистента фермера на основі LLM і економічного обґрунтування прийнятих рішень використовувався аналітичний метод досліджень, а для розробки асистента фермера на основі LLM – моделювання.

Інформаційна база кваліфікаційної роботи сформована з Інтернет-ресурсів, що містять інформацію про NLP, LLM, нейронні мережі на основі архітектури трансформера, інструментарій для локалізації LLM.

Практична значущість роботи полягає у розробці асистента фермера на основі великих мовних моделей – може бути використаний для подальших досліджень за даною тематикою та при проектуванні розумної ферми.

Апробація результатів відбувалася в рамках XIX щорічної студентської наукової конференції «Сучасні інформаційні технології та інноваційні методики в економіці, менеджменті та бізнесі» Полтавського державного аграрного університету (14 травня 2024 р., м. Полтава).

За результатами досліджень здійснено публікацію тез доповідей.

Структура кваліфікаційної роботи логічно пов'язана з завданнями досліджень і містить вступ, три розділи основної частини, висновки, список використаних джерел, додатки. Загальний обсяг пояснювальної записки кваліфікаційної роботи складає 59 сторінок формату А4. Вона містить 19 рисунків.

РОЗДІЛ 1

ОСОБЛИВОСТІ ЗАСТОСУВАННЯ ШТУЧНОГО ІНТЕЛЕКТУ В АГРОСЕКТОРІ

1.1 Аналіз сучасного стану розвитку штучного інтелекту

Штучний інтелект (ШІ) відноситься до області комп'ютерних наук, яка займається створенням програм та систем, спроможних вирішувати задачі, які потребують людського мислення. Ці задачі можуть охоплювати різні аспекти, включаючи розпізнавання мови, візуальне сприйняття, навчання та планування.

Простіше кажучи, ШІ – це змога машин вирішувати задачі, для яких зазвичай потрібен людський інтелект, як-от розпізнавання зорових образів, розпізнавання мови, ухвалення рішень і переклад мови. Інакше кажучи, системи ШІ здатні обробляти дані, вчитися на їхній основі та робити передбачення й ухвалювати рішення, ґрунтуючись на цих знаннях, що дає змогу машинам виконувати задачі з вищою точністю та швидкістю й ефективністю, ніж людина [1].

ШІ може бути поділена на декілька видів, кожен з яких має свої особливості і сфери застосування. До цих типів належать:

- машинне навчання;
- глибоке навчання;
- обробка природної мови;
- комп'ютерний зір тощо.

Як правило, ці технології використовують для того, щоб машини могли виявляти відповідності в даних, робити прогнози та приймати рішення на основі цих даних, а також взаємодіяти з людьми природними способами, такими як мова і текст.

Особисто для мене є цікавими інтелектуальні особисті помічники, відомі також як віртуальні асистенти – це програми, які можуть надавати персональну

інформацію, виконувати завдання та послуги для користувача. Вони використовують ШІ для максимальної ефективності. Ці асистенти здатні виконувати широкий спектр завдань, включаючи здійснення пошуку інформації в Інтернеті, управління музикою, пристроями «розумного будинку», обробку списків задач і електронної пошти, взаємодію з календарем і так далі. Ви можете взаємодіяти з цими помічниками за допомогою голосу або тексту.

Серед найвідоміших інтелектуальних асистентів – Siri від Apple, Alexa від Amazon, Cortana від Microsoft та Google Assistant від Google [2]. Зараз спостерігаємо активне зростання ринку віртуальних помічників: вони стають дедалі більш поширеними в операційних системах і необхідними для досягнення максимальної ефективності в повсякденних задачах. Сучасні інтелектуальні асистенти вносять значний вклад в розвиток Інтернету речей, вони спрощують взаємодію між людьми і технологіями, роблячи цю взаємодію більш природньою та ефективною. Вони перетворюють наші пристрої, автомобілі, будинки на «розумні» об'єкти, що значно поліпшує наше повсякденне життя.

Ось кілька способів, якими вони це роблять:

- часто люди використовують голосові команди для управління такими асистентами, як Siri, Alexa, Google Assistant та інші, це робить процес взаємодії з технологіями простим і природнім, так як вам не потрібно знати специфічні команди або інструкції;

- ці помічники можуть виконувати декілька завдань одночасно, створюючи ефективне робоче середовище для користувачів, наприклад, вони можуть відтворювати музику, перевіряти погоду, контролювати «розумний будинок» та навіть керувати вашим календарем одночасно;

- інтелектуальні асистенти здатні вивчати звички, вподобання та доступ до інформації користувача для надання персоналізованої інформації та рекомендацій, це може значно поліпшити взаємодію користувача з технологіями;

– інтелектуальні асистенти використовують обробку природної мови (Natural Language Processing, NLP) для розпізнавання та відповіді на людські команди, оскільки взаємодія з масштабними базами даних або командами може бути складною, NLP допомагає скоротити цей бар'єр, перекладаючи людські команди на дії, які асистент може виконати;

– велика можливість для віртуальних асистентів – це здатність інтегруватися з різними системами та платформами, що дозволяє користувачам ефективно керувати всіма своїми цифровими середовищами через одного помічника.

Історія ШІ сягає давніх цивілізацій, де в казках і міфах зображували машини та істот із людським інтелектом. Однак сучасні дослідження в галузі ШІ почалися тільки в 1950-х р., коли з'явилися перші електронні комп'ютери: 1956 р. група дослідників із Дартмутського коледжу організували семінар зі «ШІ», який загальноприйнято вважати початком цієї галузі. За десятиліття, що минули відтоді, дослідження в галузі ШІ досягли значного прогресу, включно з розробкою експертних систем, нейронних мереж та алгоритмів машинного навчання, але прогрес був нерівномірним, і дослідження в області ШІ переживали періоди наснаги і невдач. В останні роки завдяки збільшенню обчислювальної потужності та величезному обсягу даних інтерес до ШІ відродився область знову активно рухається до своєї остаточної мети - створення розумних машин, які зможуть змагатися з людським інтелектом і можливо перевершувати його.

2024 р. став роком знакових подій та відкриттів у сфері ШІ. Він вносить постійні зміни в різні сфери нашого життя, включаючи ринок праці, маркетинг, і навіть наше сприйняття реальності. Ключові тенденції у сфері ШІ 2024 р. це розвиток генеративних моделей Штучного Інтелекту. Завдяки великим анонсам від гігантів індустрії, таких як OpenAI і Microsoft, можна очікувати значне прискорення розвитку сфери [22].

Також у 2024 р. може статись прорив в області «Deep Fakes» або гіперреалістичних фейкових зображень і відео, згенерованих за допомогою

ШІ. Наприклад, картинка Папи Римського у стильній дорогій куртці від Balenciaga (рис. 1.1) стала вірусною і чимало людей повірили у її справжність. Також широкий резонанс мало зображення з фальшивим арештом Д. Трампа.



Рисунок 1.1 – Згенероване зображення Папи Римського

Це може призвести до серйозних викликів щодо вірогідності візуальних даних. Потужний вплив ШІ має важливі наслідки для ринку праці та маркетингу. З однієї сторони, ШІ може замінити деякі традиційні професії, проте з іншої сторони, він також створює нові робочі місця та можливості. В маркетинговій сфері ШІ може допомогти в аналізі тенденцій, в підборі аудиторії та в особистому маркетингу. Проте не дивлячись на всі бар'єри і виклики, які приносить штучний інтелект, він продовжує відкривати нові горизонти досліджень та можливостей, накопичуючи знання ще більш потужних і корисних програм та додатків для нас всіх [29].

А також 2023 р. був проривним роком для ШІ. Від впливу ШІ на військову сферу до його впровадження в комерційний сектор, цей період був інтенсивним стосовно нових досягнень у галузі ШІ. Поява генеративного ШІ викликала зміни в різних сферах, включаючи спам, шахрайство та дискредитацію відомих людей. Проте, випереджати прогрес не варто: з його впровадженням відкриваються нові можливості для створення більш

ефективних систем, здатних виконувати складні завдання. Українська ІТ-компанія Brainstack_ підкреслила основні тренди ШІ 2023 р., в їх числі зростаюча відповідальність і повага до етики, більш активне використання передових машин, що навчаються і звичайно ж, посилення життєздатності технології штучного інтелекту. В 2023 р. ШІ застосовувався в Україні в різних сферах від військової до медіа. Це стало свідченням глобального підходу до технології і значної ролі ШІ в передових рішеннях.

Штучний інтелект продовжує розвиватися з кожним роком. Цей розвиток включає в себе нові можливості і виклики, але вихід за обмеження, якими був навчений ШІ, створює можливості для нових відкриттів і застосувань, які мають потенціал знизити витрати, покращити продуктивність і збільшити доступність до різних служб і досліджень. Бадьорий погляд на майбутнє – це те, що заохочує нас продовжувати дослідження і розвиток ШІ.

1.2 Прикладні аспекти застосування штучного інтелекту в сільському господарстві

Сільське господарство – це одна з найважливіших галузей, яка забезпечує людство продуктами харчування. Але сучасне сільське господарство зіштовхнулося із серйозними викликами, такими як зміна клімату, недолік ресурсів та нестабільність цін на сировину. У цьому контексті штучний інтелект може висвітлити новий шлях для сільського господарства, забезпечуючи покращення ефективності та прибутковості [23].

ШІ може допомогти вирішити задачі, які раніше були складно виконати за допомогою традиційних методів. Наприклад, системи навігації із використанням ШІ можуть допомогти фермерам знаходити найкращі маршрути для поливу та розподілу ресурсів, таких як добрива та пестициди. Також, штучний інтелект може допомогти в управлінні високотехнологічною сільськогосподарською технікою, такою як дрони і автоматизовані трактори.

Наприклад система автономного керування спецтехнікою. Штучний інтелект у разі застосовується для аналізу даних, що надходять з камер та інших систем, що дозволяє просування комбайнів, тракторів або обприскувачів успішно аналізувати тип і положення кожного об'єкта. Це покращує ефективність роботи, зменшує витрати та допомагає підтримувати більш екологічно чисті методи сільськогосподарського виробництва.

Сенсори та датчики:

- камери використовуються для візуального аналізу оточення, вони можуть визначати типи об'єктів (наприклад, рослини, перешкоди) і їх положення;

- GPS-модулі забезпечують точне позиціонування техніки на полі;

- лідари та радари вимірюють відстані до об'єктів та створюють тривимірні карти оточення;

- інші сенсори можуть вимірювати вологість ґрунту, рівень поживних речовин, температуру та ін.;

Аналіз даних:

- інформація з усіх сенсорів збирається та передається на обробку;

- використання алгоритмів машинного навчання для обробки даних, ШІ аналізує зображення, визначає типи рослин, оцінює їх стан та визначає перешкоди;

- алгоритми машинного навчання постійно оновлюються та вдосконалюються на основі нових даних, що дозволяє системі ставати точнішою з часом.

Прийняття рішень:

- система приймає рішення щодо руху техніки, наприклад, уникнення перешкод, вибір оптимального маршруту для обробки поля, регулювання швидкості та інших параметрів роботи;

- система може регулювати процеси внесення добрив, води та засобів захисту рослин, забезпечуючи точність та економію ресурсів.

Переваги використання такого виду технологій:

- покращення ефективності, детальний контроль за обладнанням та оптимізація робочих процесів призводять до зниження споживання палива, добрива та інших матеріалів;

- автономні системи зменшують потребу в персоналі та знижують ймовірність помилок, зумовлених людським фактором;

- ретельне дозування добрив та засобів захисту рослин допомагає знизити їх шкідливий вплив на навколишнє середовище;

- обережний догляд за рослинами і вдосконалення умов їх вирощування сприяють покращенню продуктивності та якості врожаю.

Приклади застосування:

- автономні комбайни самостійно збирають врожай, визначаючи найефективніші маршрути та уникають потенційних пошкоджень врожаю;

- трактори використовуються для виконання різноманітних аграрних робіт, зокрема оранки, сівби та обприскування, автономні трактори забезпечують високу точність виконання цих операцій;

- обприскувачі це пристрої для точного обприскування рослин зменшують використання хімічних реагентів і сприяють збереженню чистоти екосистеми.



Рисунок 1.2 – Сучасні технології сільського господарства

Одним з основних викликів сучасного сільського господарства є погода. ШІ може допомогти у передбаченні погодних змін з високою точністю та розробити стратегії для максимально ефективного використання ресурсів.

Наприклад, ШІ може проводити аналіз даних щодо ґрунту та вологості, щоб допомогти фермерам визначити найкращий час для посіву та збору урожаю [24].

Іншим важливим аспектом сільського господарства є захист врожаю від шкідників та хвороб. Системи із штучним інтелектом можуть розпізнавати зайву рослинність та шкідників, а також визначати оптимальний час для застосування пестицидів. Це допомагає зберегти ресурси та зменшити вплив на навколишнє середовище.

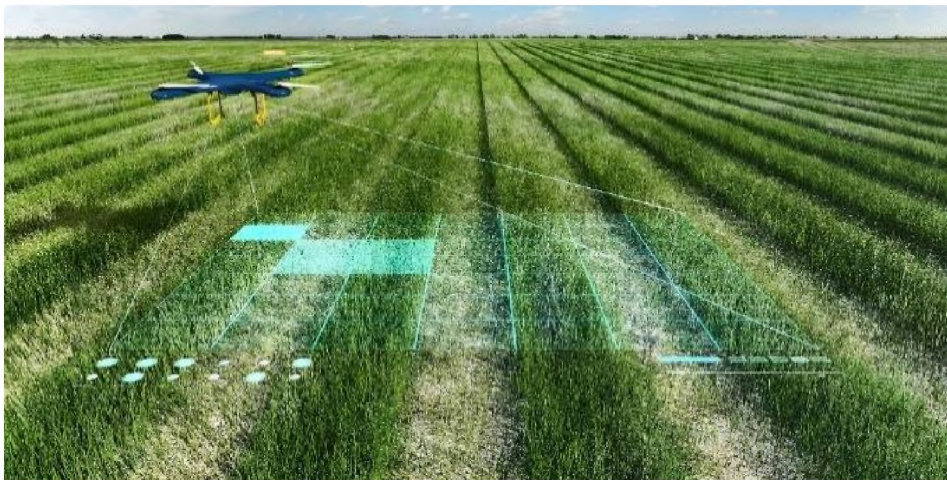


Рисунок 1.3 – Прогнозування та аналітика за допомогою ШІ

Нарешті, ШІ може використовуватися для передбачення попиту на продукти та ринкової ціни. Це допомагає фермерам збалансувати свої посівні площі та забезпечити ефективне використання ресурсів.

В цілому, він може бути потужним інструментом для досягнення стійкого та ефективного сільського господарства. Величезна кількість даних, що збираються в реальному часі, здатні бути оброблені і аналізовані штучним інтелектом для прийняття швидких та точних рішень, що допомагають фермерам підтримувати здоров'я своїх рослин та зберігати навколишнє середовище. ШІ має великий потенціал для розвитку сільського господарства, який допоможе забезпечити сталість та продуктивність для наступних поколінь [3].

Однак, виникає питання етичності використання ШІ у сільському господарстві. Наприклад, чи може штучний інтелект замінити людей в розробці технологій для сільського господарства? Чи може він змінити традиційний спосіб життя селян та поглинути ринок праці? Одним з рішень цього питання може бути розвиток програм та тренування фермерів у використанні штучного інтелекту для поліпшення своїх процесів. Також важливо забезпечити доступність технологій ШІ для всіх фермерів, незалежно від їх географічного розташування чи економічного статусу.

У певних областях, інтегрування штучного ШІ в сільське господарство може ще більше сприяти ефективності та екологічності виробництва. Наприклад, точне землеробство, яке використовує дані з дронів та супутників для аналізу поля, дозволяє оптимізувати використання води та добрив, мінімізуючи втрати та шкідливий вплив на довкілля. Це демонструє, яким чином ШІ може сприяти розв'язанню всесвітніх викликів, такі як продовольча безпека та зміна клімату, пропонуючи інноваційні рішення для підвищення продуктивності та стійкості агросектору.

Штучний інтелект можна з успіхом використовувати для прогнозування врожаю в різних регіонах, незалежно від специфіки їх клімату, ґрунтів чи особливостей культур. Адже основою для тренування алгоритмів штучного інтелекту є великі обсяги інформації, зібраної з різноманітних ресурсів та в усіх точках світу.

Наприклад, ШІ може використовувати різні дані, включаючи історичні дані про врожай, дані з супутникових зображень, дані про температуру і вологість повітря, типи ґрунтів, інформацію про висаджені культури та інше, для створення точних прогнозів врожаю для різних регіонів.

Варто також наголосити, що системи, засновані на ШІ постійно «навчаються», тому їх прогнози стають все більш точними з часом.

Але в їхніх прогнозах також може бути враховано введення нових технологій сільського господарства, керування водними ресурсами задля іригації, погодні умови та навіть соціально-економічні фактори. Тому,

незважаючи на регіональні особливості, ШІ може стати потужним інструментом для прогнозування врожаю в будь-якому регіоні світу.

У підсумку, штучний інтелект має великий потенціал для розвитку сільського господарства у майбутньому. Використання ШІ може допомогти фермерам збільшити ефективність та прибутковість, мінімізувати вплив на навколишнє середовище та забезпечити сталу продукцію для населення. Проте, важливо забезпечити етичне та доступне використання штучного інтелекту, а також розвивати нові методи та програми для навчання фермерів його використанню. Тільки тоді ШІ може стати справжньою революцією для сільського господарства, покращуючи життя людей та захищаючи наше навколишнє середовище.

1.3 Реалізація штучного інтелекту на основі нейронних мереж з архітектурою трансформер

Трансформер є типом архітектури нейронної мережі, створеної для задач NLP. Останнім часом він став однією з найбільш розповсюджених архітектур у цій сфері [13].

Ця архітектура базується на концепції само-уваги (self-attention), яка дозволяє моделі концентрувати увагу на різних сегментах вхідної послідовності під час обчислення представлення кожного слова чи лексеми. Ідея само-уваги запозичена з людського способу розуміння об'єктів або тексту шляхом зосередження лише на певних частинах. Наприклад, замість звертання уваги на всі деталі зображення, людина зосереджується лише на конкретних його частинах для кращого розуміння. Подібним чином модель може фокусуватися лише на найважливішій інформації для досягнення кращого розуміння.

У традиційних архітектурах нейронних мереж, таких як рекурентні нейронні мережі (RNN), вхідна послідовність обробляється послідовно, слово

за словом. Трансформер, навпаки, може обробляти всю вхідну послідовність паралельно, що робить його набагато ефективнішим для довгих послідовностей.

Попередньо навчені мовні моделі (Pre-Trained Language Models, PTLM) навчають трансформери на великих корпусах текстів, передаючи отримані знання для таких задач, як автоматичне підсумовування текстів (ATS). Завдяки цьому багаті семантичні та контекстуальні характеристики вбудованих слів, отриманих за допомогою PTLM, покращують якість кінцевих анотацій. Широке застосування трансформерів та PTLM у різних задачах NLP, включаючи абстрактне автоматичне підсумовування, робить їх основою сучасних досліджень у NLP [4].

Трансформер складається з набору енкодерів і декодерів. Кожен енкодер містить два шари: шар само-уваги і шар нейронної мережі прямого поширення. У декодері присутні ті ж самі шари, але між ними додається шар уваги, який дозволяє декодеру фокусуватися на відповідних частинах вхідного речення. На рисунку нижче представлена схема архітектури трансформера (рис. 1.4).

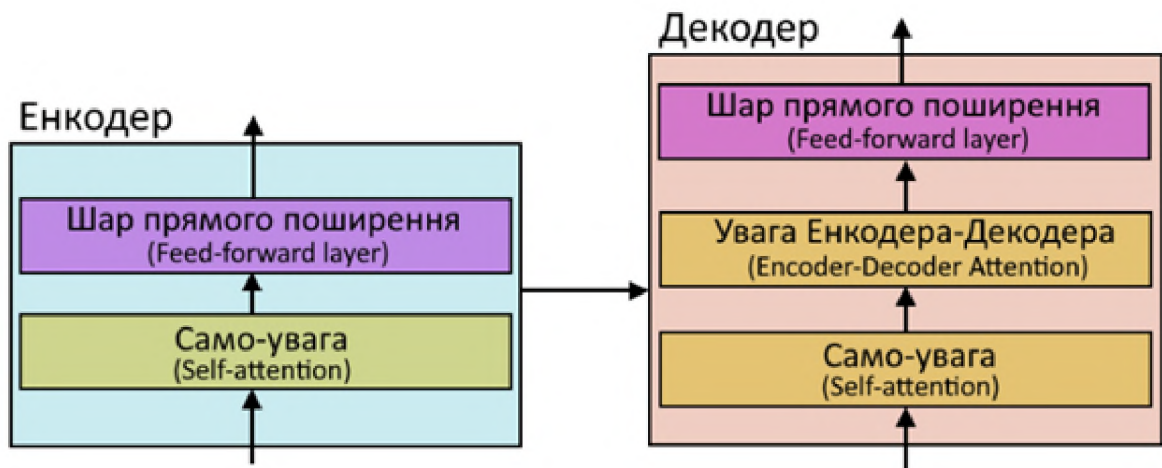


Рисунок 1.4 – Схема архітектури трансформера

Щоб мовна модель могла передбачити значення тексту, вона повинна враховувати контекстуальну подібність слів. Речення, що надходять до трансформера, перетворюються на вбудовування слів – низьковимірні векторні представлення тексту, які зберігають контекстну подібність слів. Ці вбудовування слів обробляються за допомогою навчальних алгоритмів, що робить їх ефективнішими та виразнішими представленнями слів.

Процес вбудовування відбувається лише в самому нижньому енкодері. Усі енкодери отримують список векторів розміром 512. Для нижнього енкодера це вбудовування слів, а для інших енкодерів — вихід попереднього енкодера. Розмір цього списку векторів є гіперпараметром, тобто параметром, що керує процесом навчання та визначає значення параметрів моделі, які вивчає алгоритм навчання.

Кожне слово у вхідному реченні проходить через процес само-уваги, після чого кожне з них обробляється окремо за допомогою прямого нейронного зв'язку (рис. 1.5). Однією з основних переваг трансформера є можливість паралельної обробки слів, що значно підвищує ефективність.

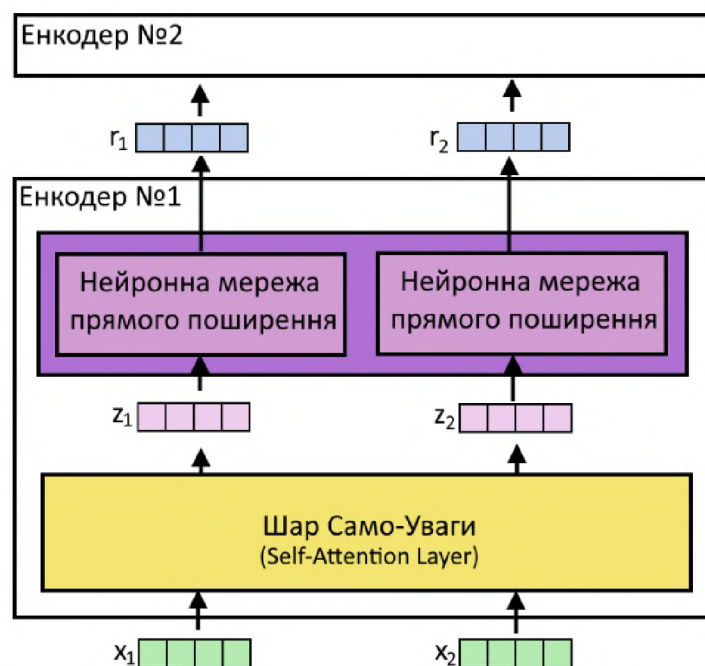


Рисунок 1.5 – Схема обробки слів енкодерами

Спосіб, яким само-увага працює, полягає в тому, що модель може розуміти, коли звертаємось до об'єктів неявно. Модель аналізує кожне слово в вхідній послідовності і може оглядати інші слова в цій послідовності, щоб отримати підказки, які допоможуть з кращим кодуванням цього слова.

Шар само-уваги приймає вхідний текст і створює три вектори для кожного слова, що називаються запитом, ключем і значенням. Ці вектори формуються шляхом множення вбудовування на три матриці, які визначаються під час навчання моделі. Розмірність вхідних векторів – 512, а обчислені вектори – 64. Це зроблено навмисно, щоб зробити обчислення багатосторонньої уваги більш ефективними. Оцінка кожного слова в порівнянні з іншими словами в реченні обчислюється шляхом скалярного добутку вектора запиту на вектор ключа відповідного слова. Ці оцінки діляться на 8 для отримання стабільніших градієнтів, і потім відбувається нормалізація за допомогою функції softmax, щоб всі вони були позитивними та склалися в 1. Кожен вектор значення множиться на відповідну оцінку softmax. Таким чином, зберігаються значення слів, на яких модель зосереджується, і заглушуються нерелевантні слова. Вихід з рівня само-уваги для кожного слова є сумою всіх зважених векторів значення. Для прискорення обробки матриці використовуються замість окремих векторів.

Таким чином, розглянули сутність та особливості архітектури трансформера, механізм само-уваги є ключовими елементами моделі, що розробляється. Вони є важливими складовими для розуміння та реалізації ШІ на основі нейронних мереж з архітектурою трансформера.

Надалі, доцільно застосовувати їх використовувати в якості основи для створення персонального асистента фермера.

РОЗДІЛ 2

ДОСЛІДЖЕННЯ ІНСТРУМЕНТАРІЮ ДЛЯ СТВОРЕННЯ ПЕРСОНАЛЬНОГО АСИСТЕНТА ФЕРМЕРА НА ОСНОВІ ВЕЛИКИХ МОВНИХ МОДЕЛЕЙ

2.1 Визначення великих мовних моделей

Початки великих мовних моделей (Large Language Model), глибокого навчання та машинного перекладу датуються стародавніми міфами та легендами про створення штучних істот, яких майстри обдаровували інтелектом або свідомістю. Сучасний ШІ з'явився у результаті майстерності філософів, котрі намагалися описати процес людського мислення як механічну маніпуляцію символами [5].

Глибоке навчання, підмножина методів машинного навчання, заснована на штучних нейронних мережах з навчанням подань.

Основоположники штучного інтелекту такі, як Ф. Розенблатт, Д. Маккарті та М. Мінській, зробили величезний внесок у розвиток цієї області. Наприклад, Розенблатт створив Перцептрон – один з перших алгоритмів машинного навчання, а Маккарті та Мінський були серед засновників області штучного інтелекту і внесли величезний внесок у її розвиток.

А. Т'юрінг також відіграв важливу роль у розвитку машинного інтелекту, будучи одним з перших, хто здійснив істотні дослідження в цій галузі. Щодо Р. Гартмаєра та його впливу на ШІ, інформація може бути обмеженою, і її слід шукати додатково.

Загалом, ці піонери штучного інтелекту були свого роду магами своїх часів, використовуючи свою експертизу та творче мислення для розвитку нашого розуміння і використання ШІ у сучасному світі. Вивчення їхніх робіт до сих пір слугує основою для розвитку нових методів в галузі штучного інтелекту та машинного навчання. Від початків до сьогодні, еволюція

глибокого навчання представляє собою послідовність важливих етапів та інновацій, які сприяли зацікавленості в цій галузі та значно розширили її можливості. Одним з критичних відкриттів в розвитку LLM була поява глибоких нейронних мереж. Вперше ідея використання мереж з багатьма шарами з'явилася в 1960-х, але до 2000-х років вона була мало використана через обмеженість обчислювальних ресурсів та недостатність великих наборів даних. Алгоритми згорткових нейронних мереж (Convolutional Neural Network) зайняли важливе місце в розвитку LLM у сфері комп'ютерного зору, дозволяючи досягти значних успіхів у завданнях розпізнавання образів та класифікації. Рекурентні нейронні мережі (Recurrent Neural Networks) також внесли свій внесок у розвиток LLM, особливо в області обробки природної мови та машинного перекладу, завдяки їхній здатності моделювати послідовні дані та враховувати контекст [28].

Крім цього, важливими досягненнями в розвитку LLM були методи навчання, такі як зворотне розповсюдження помилки та методи оптимізації, які сприяли ефективному навчанню глибоких нейронних мереж на великих обсягах даних і покращили їхню точність та швидкодію. Загалом, завдяки появі глибоких нейронних мереж, алгоритмів CNN та RNN, а також вдосконаленню методів навчання, LLM зазнала значних досягнень і стала основою для багатьох передових застосувань у сферах розпізнавання образів, обробки природної мови, машинного перекладу та інших.

Обчислювальні системи, відомі як штучні нейронні мережі, були створені під впливом біологічних нейронних мереж, що формують мозок тварин. Такі мережі виникли з набору сполучених вузлів, або штучних нейронів, які імітують нейрони в біологічному мозку. Всі з'єднання, подібні до синапс у природного мозку, можуть переслати сигнал до інших нейронів. Штучний нейрон приймає сигнали, переробляє їх і може посилати сигнали до нейронів, з якими він з'єднаний.

Складніші нейронні архітектури, зокрема глибокі нейронні мережі, були розроблені для розв'язання більш складних завдань. Останніми роками

великий інтерес викликав глибоке навчання – підмножина методів машинного навчання, заснована на штучних нейронних мережах з навчанням подань. Глибоке навчання включає використання глибоких нейронних мереж, з'єднання багатьох штучних нейронів у багат шарову структуру [9].

Структурні зміни в архітектурі нейронних мереж дозволили штучному інтелекту робити прогнози та приймати рішення з високою точністю і ефективністю. Завдяки цьому технології, такі як машинний переклад, обробка природної мови та комп'ютерне зору, досягли нових вершин. Багато видів нейронних мереж, включаючи згорткові та рекурентні нейронні мережі, було створено для конкретних завдань машинного навчання. CNN, що імітують зоровий кортекс, виявилися особливо ефективними для задач, пов'язаних із зображеннями та відео. RNN, які можуть обробляти дані послідовно, дозволяють виконувати завдання подібні до генерації тексту або машинного перекладу.

Зворотне розповсюдження помилки – це класичний метод навчання, який використовує градієнтний спуск для оптимізації ваг у мережі на основі помилки між очікуваним та дійсним виведенням. Це основний метод навчання для більшості сучасних нейронних мереж.

Розгляд ключових досліджень та наукових праць, які мають вагомий внесок у розвиток глибокого навчання (LLM), охоплює аналіз визначних вчених, методологій та відкриттів, що суттєво вплинули на подальший прогрес у цій галузі.

Одним із важливих досліджень є розробка «Методу зворотного розповсюдження помилки», запропонованого Д Рамелхартом, Д Румельхартом та Т Сімоном в 1986 р. Цей підхід визначив нові стандарти для навчання глибоких нейронних мереж і відкрив шлях для їхнього широкого використання. Іншим вагомим досягненням стало введення моделі «Deep Belief Networks» (DBN) Геоффри Гінтоном та колегами у 2006 р. Це дослідження відкрило нові перспективи у навчанні глибоких нейронних мереж та дало поштовх для розвитку різних напрямків досліджень. Також варто

зазначити роботу Я. Лекуна, Г. Хінтона та К. Сундерхаузера з розробки «Конволюційних нейронних мереж» в 1998 р., які застосовуються у сфері аналізу візуальних зображень. Ця праця стала основою для подальшого розвитку CNN у різних областях, зокрема у комп'ютерному зорі та обробці зображень. До того ж, слід відзначити внесок Й. Бенджіо та його колег, які працювали над теоретичними та алгоритмічними інноваціями, що допомогли розкрити потенціал глибокого навчання та покращити його результати в різних задачах.

Ці ключові наукові дослідження та роботи відображають лише деякі з ряду важливих внесків у галузі глибокого навчання, свідчаючи про постійний прогрес та створення основ для подальшого розвитку цієї динамічної науки.

LLM стали визначальним етапом в розвитку сучасного штучного інтелекту та машинного навчання. Це алгоритм глибокого навчання, який може розпізнавати, узагальнювати, перекладати, прогнозувати та генерувати тексти та інший вміст на основі обробки природної мови [14].

Перший крок до створення LLM включає підготовку даних для тренування моделі. Великі мовні моделі, як правило, навчаються на великих наборах текстових даних. Структура або архітектура моделі визначає, як дані обробляються в мережі.

Другим важливим аспектом роботи з LLM є оптимізація моделі. Це включає процеси, які направлені на вдосконалення продуктивності та ефективності моделі, враховуючи її точність та швидкість. Оцінка дозволяє перевірити, наскільки добре модель виконує завдання, для якого вона була створена.

LLM – це мовні моделі, що працюють з великою кількістю додаткових параметрів. Це означає, що вони використовують велику кількість даних та мереж, що дозволяє їм створювати більш точні та своєрідні висновки (моделі LLM та можливості, які вони відкривають) [14].

Останнім, але не менш важливим, є питання етики при роботі з великими мовними моделями. Важливо розуміти, як використовувати цей технологічний

прогрес відповідально, враховуючи можливі наслідки (8 етичних міркувань LLM) [20].

Якщо деталізувати дослідження, може знадобитися консультація вчених та дослідників, що працюють в даній галузі, оскільки їхня робота може містити конкретні методи і відкриття, які надають значний вклад у розвиток LLM.

2.2 Оцінка властивостей великих мовних моделей

Принципи функціонування та основні функції технології можливо представити кількома положеннями.

Основний принцип функціонування LLM полягає в використанні нейронних мереж, які спрямовані на моделювання роботи людського мозку. Вони складаються зі штучних нейронів, організованих у великій кількості шарів. Ця глибока архітектура дозволяє системі автоматично відбирати та аналізувати корисні ознаки з вхідних даних на різних рівнях абстракції, що сприяє здійсненню точних прогнозів та класифікації об'єктів.

Поширення помилок через навчання : Процес навчання LLM ґрунтується на алгоритмі зворотного поширення помилки. Під час цього процесу, система порівнює прогнозовані значення зі справжніми, і коригує ваги нейронів так, щоб мінімізувати помилку. Це дозволяє системі підганяти свої параметри для кращого відображення вхідних даних та здійснення точних передбачень.

Автоматизація відбору ознак : LLM володіє унікальною здатністю автоматично відбирати корисні ознаки з вхідних даних, що відіграє важливу роль у здійсненні різноманітних завдань, таких як розпізнавання образів, розпізнавання мови, аналіз тексту та інші.

Постійне навчання і адаптація: Ще одним важливим принципом функціонування LLM є його здатність до постійного навчання та адаптації до змінних умов. Це означає, що система може навчатися на нових даних та

вдосконалювати свої навички з часом, що робить її ефективнішою та динамічнішою в різних сценаріях застосування [17].

Процес освоєння LLM включає декілька важливих етапів і може бути викликом. Ось спрощений покроковий план цього процесу:

Використовуючи велику кількість текстових даних, які можуть бути згруповані з джерел, таких як книги, веб-сторінки, статті або соціальні медіа, починається навчання LLM. Мета – охопити багату різноманітність людської мови.

Очищення даних: необроблені текстові дані впорядковуються в процесі, який називається попередньою обробкою. Це включає такі завдання, як видалення непотрібних символів, розбиття тексту на невеликі фрагменти, які називаються маркерами, і перетворення всього цього у формат, який може використовуватися моделлю.

Розділення даних: чисті дані поділяються на 2 набори. 1. Один набір навчальних даних використовується для навчання моделі. Інший набір, дані перевірки, пізніше використовується для тестування продуктивності моделі.

Налаштування моделі: визначається структура LLM, яка називається архітектурою. Це включає вибір типу нейронної мережі та прийняття рішень щодо різних параметрів, таких як кількість шарів у мережі та прихованих елементів.

Навчання моделі: тепер починається ваше власне навчання. Модель LLM переглядає та засвоює навчальні дані, робить прогнози на основі вивченого на даний момент та коригує внутрішні параметри, щоб зменшити різницю між прогнозом та фактичними даними.

Перевірка моделі: навчання моделі LLM перевіряється за допомогою даних перевірки. Це дозволяє побачити, наскільки добре працює ваша модель, і налаштувати параметри моделі для підвищення продуктивності.

Використання моделі: після навчання та оцінки ви готові до використання моделі LLM. Тепер ви можете інтегрувати його в додаток або систему, яка генерує текст на основі нещодавно введених даних.

Покращення моделі: нарешті, завжди є можливості для вдосконалення. З часом моделі LLM можна вдосконалювати, використовуючи оновлені дані або коригуючи налаштування на основі відгуків та фактичного використання.

Цей процес вимагає критичних обчислювальних ресурсів, таких як потужні процесори та великий обсяг пам'яті, а також професійні знання машинного навчання. З цієї причини це зазвичай роблять спеціалізовані науково-дослідні установи або компанії, які мають доступ до необхідної інфраструктури та досвіду.

Створений OpenAI, GPT-3 був навчений на величезній кількості даних, включаючи текст з різних джерел. Зокрема, модель GPT-3 була навчена на 570 GB текстової інформації, але точний обсяг набору даних не вказано. В основу цього тексту лягли різні джерела: книги, статті, веб-сторінки, енциклопедії, форуми та ін. Важливо відзначити, що в GPT-3 використовується суміш ліцензійних даних, опублікованих даних і даних, які дозволено використовувати OpenAI [25].

Навчання цієї моделі проводилося з використанням дуже потужних обчислювальних ресурсів і зайняло багато часу. Така велика кількість даних дозволяє моделі розуміти та генерувати текст у різних контекстах, що робить її дуже універсальною та здатною виконувати складні завдання обробки природної мови. Меншим LLM може знадобитися менше – можливо, 10-20 GB або навіть 1 GB – але це все одно багато.

В рамках наукового дослідження, важливо наголосити, що LLM вже не є просто абстрактною теорією або експериментальним концептом. Все більше й більше вони виступають як фундаментальний елемент у нашому цифровому контексті. Розглянемо деякі визначальні характеристики.

Майстерність імітації людського тексту. LLM перетворюють засоби розробки мовних задач. Розроблені за допомогою потужних алгоритмів машинного навчання, ці моделі вирізняються можливістю розпізнавати та інтерпретувати нюанси людської мови, включаючи контекст, емоціональні відтінки та, у певній ступені, сарказм. Ця спроможність відтворювати манеру

людського спілкування не просто новаторський бенчмарк, вона надзвичайно важлива. Просунуті письмові здібності LLM можуть удосконалити численні аспекти, від розробки контенту до обслуговування клієнтів.

Уявіть, що ви маєте можливість поставити складне питання вашому цифровому асистенту і отримати відповідь, що не просто зрозуміла, але й послідовна, відповідна контексту і подана природнім для людини способом. Ось що нам дають LLM. Вони спонукають до більш інтуїтивно зрозумілої й привабливої взаємодії між людиною та комп'ютером, посилюють користувацький досвід і демократизують доступ до важливих даних.

Доступна обчислювальна потужність. Прогрес LLM не міг би відбутися без одночасних досягнень в області комп'ютерних технологій. Насамперед, широкий доступ до високоефективних обчислювальних ресурсів відіграв ключову роль у їхній еволюції та прийнятті [6].

Сучасні хмарні платформи надають неперевершений доступ до потужних обчислювальних ресурсів. Таким чином, дослідники та малі організації мають можливість тренувати складні моделі машинного навчання.

Новітні винайдення в області обчислювальних процесорів (наприклад, GPU та TPU) та розвиток розподілених обчислень сприяли навчанню моделей з мільярдами параметрів. Зростання доступності обчислювальних ресурсів стимулює розвиток та успіх LLM, що призводить до постійного поступу та розширення застосувань у цьому напрямку.

Зміна споживчих уподобань. Сучасні споживачі не просто прагнуть отримати відповіді; вони хочуть привабливих та персоналізованих взаємодій. Через те, що все більше людей виростає з використанням цифрових технологій, стає очевидною зростаюча потреба в технологіях, що виявляються більш натуральними та людськими. LLM надає надзвичайну можливість задовольнити ці очікування. Генеруючи текст, що відчувається як написаний людиною, ці моделі можуть формувати захоплюючі та жваві цифрові враження, які можуть сприяти задоволеності та лояльності користувачів. Чи то штучно інтелектуальні чат-боти допомагають клієнтам, чи голосові асистенти

оновлюють новинні дайджести, LLM відкривають епоху ШІ, який ліпше пізнає нас.

Сховище неструктурованої інформації. Неструктурована інформація, наприклад електронна пошта, пости в соціальних мережах та відгуки від клієнтів, є справжнім джерелом знань. Припускається, що близько 80 % корпоративних даних є неструктурованими, при цьому їх обсяг збільшується на 55 % щорічно. Ці дані можуть стати справжньою золотоносною жилотою для компанії, якщо їх правильно використовувати. Тут починає виявлятися потенціал LLM, які спроможні обробляти та інтерпретувати такі дані у великому масштабі. Вони в силах робити різноманітні речі, такі як аналізувати емоції, класифікувати текст та витягувати дані, надаючи таким способом цінні відомості. Неважливо, чи стосується це визначення тенденцій у соцмережах чи оцінки настроїв клієнтів через відгуки - LLM допомагають бізнесам орієнтуватися в океані неструктурованих даних та приймати обґрунтовані рішення.

Розширення ринку NLP (рис. 2.1). Потенціал LLM відображається у швидко зростаючому ринку NLP. Аналітики прогнозують розширення ринку NLP 11 мільярдів доларів у 2020 р. до понад 35 мільярдів доларів до 2026 р. [26].

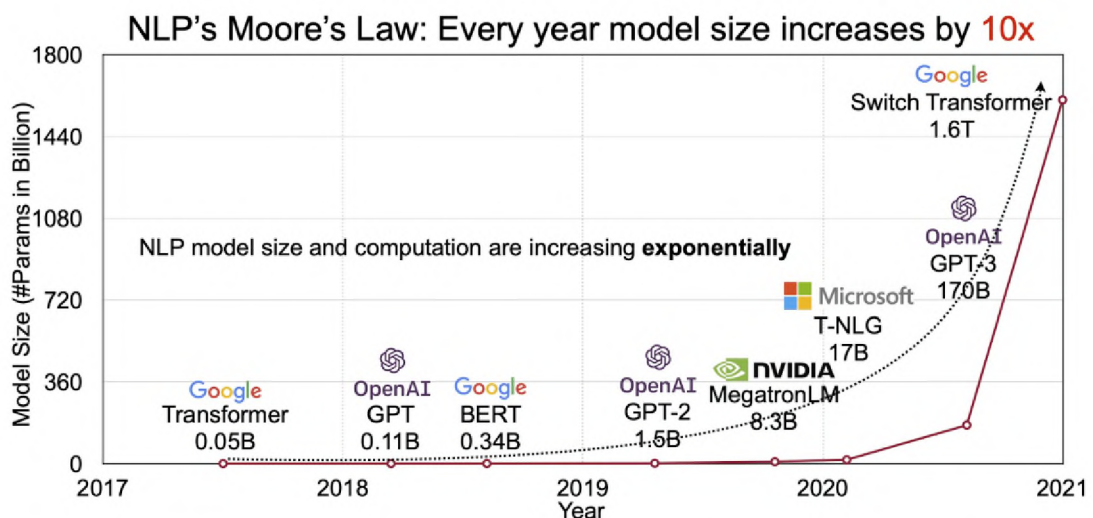


Рисунок 2.1 – Відповідно до закону Мура, розмір моделі щорічно збільшується в 10 разів.

Не тільки ринок продовжує розширюватися, але й самі моделі поступово збільшуються, як у фізичних розмірах, так і в кількості параметрів, що їх оброблюють. Тенденції росту і розвитку LLM відображені на графіку, який підкреслює їх зростаючу складність і потужність.

В майбутньому великі мовні моделі обіцяють принести захоплюючі досягнення і вчені зможуть робити дослідницькі прориви, які дозволять розширити можливості та застосування систем ШІ.

Дослідники у сфері ШІ зосереджують свої зусилля на таких важливих аспектах, як оптимізація ефективності моделі, мультимодальне навчання, персоналізація, етичний ШІ та надійність, продовжуючи розширювати межі досяжного і відкриваючи шлях до нової епохи інновацій на основі штучного інтелекту, які приносять користь як окремим користувачам, так і всьому суспільству.

2.3 Особливості застосування великих мовних моделей

На даний час, існують різноманітні приклади застосування LLM у різних галузях, які доцільно проаналізувати. Це стосується від NLP до охорони здоров'я та фінансів, LLM знаходять широке застосування завдяки своїй здатності аналізувати та генерувати текстову інформацію з високою точністю. Далі наведено ключові приклади використання LLM у різних сферах.

В рамках вирішення завдань NLP мова йде про машинний переклад: LLM використовуються для автоматичного перекладу текстів з однієї мови на іншу. Сучасні моделі, такі як GPT-4o, забезпечують високу точність перекладу, зберігаючи контекст та стилістичні особливості оригіналу. LLM здатні розпізнавати та генерувати тексти на основі заданих параметрів. Це використовується для автоматичного створення статей, новин, резюме та інших текстових матеріалів. Моделі глибокого навчання можуть аналізувати

тексти з соціальних мереж (аналіз настроїв та емоцій), відгуки клієнтів та інші джерела, щоб визначити настрої та емоції, виражені у тексті. Це допомагає компаніям зрозуміти сприйняття їхньої продукції або послуг [7].

Для охорони здоров'я слід виділити теж кілька завдань. В медичній сфері LLM використовуються для діагностики захворювань і адаптації підходів до лікування, в фінансовій області – для аналізу ринкових показників і ризик-менеджменту, а в аграрному секторі – для контролю за врожаєм та ефективного використання ресурсів. Відносно цього, активне використання LLM зауважується у виробничих галузях, освіті, маркетингу, HR-управлінні та інших секторах, де вони стимулюють інновації і підвищують продуктивність. Аналіз медичних записів. LLM використовуються для автоматичного аналізу медичних записів, що допомагає лікарям швидко знаходити необхідну інформацію, діагностувати захворювання та приймати обґрунтовані рішення щодо лікування. Підтримка діагностики – LLM здатні аналізувати медичні зображення (наприклад, рентгенівські знімки або МРТ) для виявлення аномалій та допомагати лікарям у діагностиці захворювань;

Персоналізована медицина. LLM можуть аналізувати генетичні дані пацієнтів для визначення індивідуальних ризиків захворювань та розробки персоналізованих планів лікування.

Фінанси та банківська справа. Моделі глибокого навчання використовуються для аналізу фінансових ринків, прогнозування цін на акції, оцінки ризиків та прийняття інвестиційних рішень. LLM допомагають виявляти підозрілі транзакції та попереджувати шахрайство, аналізуючи великі обсяги фінансових даних у режимі реального часу. Боти на основі LLM можуть відповідати на запити клієнтів, надавати фінансові консультації та допомагати у вирішенні проблем.

Роздрібна торгівля та електронна комерція. Персоналізовані рекомендації – LLM аналізують поведінку покупців на веб-сайтах та в додатках, щоб надавати персоналізовані рекомендації щодо товарів, які можуть зацікавити користувачів. Аналіз відгуків клієнтів – моделі можуть автоматично

аналізувати відгуки про продукти та послуги, визначаючи ключові проблеми та позитивні моменти, що допомагає компаніям покращувати свою продукцію. Оптимізація ланцюга постачання – LLM використовуються для прогнозування попиту на товари, оптимізації запасів та управління ланцюгом постачання, що зменшує витрати та підвищує ефективність.

Освіта та навчання. Інтелектуальні репетитори – персоналізовані освітні платформи на основі LLM можуть адаптувати навчальні матеріали до індивідуальних потреб учнів, надаючи рекомендації та допомогу в реальному часі. Аналіз освітніх даних – LLM допомагають аналізувати великі обсяги освітніх даних для визначення трендів, оцінки успішності учнів та покращення навчальних програм. Автоматична оцінка робіт – моделі можуть автоматично оцінювати письмові роботи студентів, забезпечуючи об'єктивність та зменшуючи навантаження на викладачів.

Юридичні послуги. Аналіз юридичних документів: LLM можуть автоматично аналізувати великі обсяги юридичних текстів, виявляючи ключові положення, потенційні ризики та невідповідності. Це значно пришвидшує процес підготовки та перевірки документів. Юридичні чат-боти: Асистенти на основі LLM можуть надавати базові юридичні консультації, допомагати користувачам знаходити потрібну інформацію та заповнювати юридичні форми.

Маркетинг та реклама. Створення контенту: LLM можуть генерувати тексти для маркетингових кампаній, включаючи рекламні оголошення, пости в соціальних мережах, блоги та новини, адаптуючи стиль під цільову аудиторію. Аналіз ринкових трендів: Моделі глибокого навчання здатні аналізувати споживчі тренди, поведінку клієнтів та відгуки, що допомагає компаніям краще розуміти потреби ринку і створювати ефективні маркетингові стратегії.

Ігрова індустрія. Генерація контенту – LLM можуть створювати сюжетні лінії, діалоги для персонажів та інші елементи гри, що робить процес розробки більш ефективним. Аналіз поведінки гравців – LLM аналізують дані про

поведінку гравців для оптимізації ігрового процесу, виявлення проблем та підвищення задоволеності користувачів.

Транспорт та логістика. LLM використовуються для аналізу даних про трафік і погодні умови, що дозволяє оптимізувати маршрути транспортних засобів і знижувати витрати на паливо та час доставки. Моделі можуть передбачати попит на транспортні послуги та товари, що допомагає компаніям краще планувати свої операції та зменшувати витрати на зберігання.

Енергетика. LLM аналізують дані про споживання енергії та погодні умови для оптимізації роботи енергомереж та зниження витрат на енергію. LLM використовуються для прогнозування попиту на енергію, виявлення аномалій в роботі обладнання та планування технічного обслуговування.

Державне управління. LLM використовуються для аналізу великих обсягів законодавчих документів, що допомагає виявляти ключові питання, тенденції та прогалини в законодавстві. Асистенти на основі LLM можуть автоматично відповідати на запити громадян, надавати інформацію про державні послуги та допомагати в заповненні форм.

Електронна комерція. LLM аналізують дані про поведінку покупців для надання персоналізованих рекомендацій щодо товарів та послуг, підвищуючи ймовірність купівлі. Автоматизація обробки замовлень, включаючи відповіді на запитання покупців, відстеження доставки та вирішення проблем, що виникають.

Туризм. Асистенти на основі LLM можуть допомагати користувачам планувати подорожі, пропонуючи маршрути, готелі та розваги на основі їхніх уподобань. Чат-боти можуть відповідати на запити гостей, надавати інформацію про послуги та допомагати у вирішенні проблем під час перебування у готелі.

Автомобільна індустрія. LLM використовуються для обробки даних з сенсорів автомобіля та прийняття рішень у режимі реального часу, забезпечуючи безпечне та ефективне водіння. Системи допомоги водію (Advanced Driver-Assistance Systems, ADAS) – моделі допомагають виявляти

небезпеки на дорозі, попереджати про можливі зіткнення та надавати інші функції для підвищення безпеки водіння.

Медіа та розваги. LLM можуть створювати музичні композиції, малюнки та інші форми мистецтва, використовуючи задані стилі та параметри. Моделі можуть допомагати сценаристам генерувати ідеї для сюжетів, діалоги та інші елементи сценаріїв. LLM використовуються для написання статей, сценаріїв, новин та іншого контенту. Вони також можуть генерувати креативні ідеї для шоу, фільмів та відеоігор. Моделі можуть аналізувати дані про вподобання аудиторії, допомагаючи медіакомпаніям створювати контент, який буде найцікавішим для їхніх глядачів. LLM можуть створювати музичні треки, писати тексти пісень та навіть генерувати відеоконтент, що відкриває нові можливості для артистів та продюсерів. Використання LLM для створення інтерактивних сюжетних ліній у відеоіграх, що адаптуються до дій гравців, роблячи ігровий процес більш захоплюючим.

Наука і дослідження. LLM допомагають обробляти великі обсяги наукових даних, виявляти закономірності та робити прогнози, що сприяє відкриттям у різних галузях науки. Моделі можуть допомагати дослідникам писати наукові статті, збирати та аналізувати бібліографічну інформацію, що підвищує ефективність наукової роботи.

Агропромисловий комплекс. LLM можуть аналізувати дані з дронів та сенсорів для моніторингу стану врожаю, виявлення хвороб рослин та прогнозування врожайності. Моделі допомагають оптимізувати використання води, добрив та інших ресурсів, що сприяє підвищенню ефективності фермерських господарств.

Сфера нерухомості. LLM аналізують тенденції на ринку нерухомості, прогнозують зміни цін та допомагають у прийнятті рішень щодо купівлі або продажу нерухомості. Автоматизація обслуговування клієнтів – чат-боти можуть відповідати на запити потенційних покупців або орендарів, надавати інформацію про об'єкти та допомагати у вирішенні питань, пов'язаних з орендою або купівлею.

LLM з успіхом впроваджуються у широкому спектрі сфер, підтверджуючи свою велику практичну вартість і ефективність. Їх використання приносить численні переваги, серед яких автоматизація рутинних завдань, покращення точності прогнозів та збільшення рівня обслуговування.

Можливості для використання LLM продовжують зростати, основними каталізаторами якого є вдосконалення алгоритмів та збільшення доступних наборів даних. Це відкриває широкі перспективи для подальшої розробки та інтеграції LLM в різні сфери нашого життя, що в свою чергу сприяє прогресу та сталому розвитку суспільства.

Загалом, використання LLM у численних галузях підкреслює значний потенціал цих технологій у трансформації підходів до ведення бізнесу, проведення наукових досліджень і надання послуг, відкриваючи нові перспективи для поліпшення як робочого складу, так і життя загалом.

2.4 Еволюція великих мовних моделей для локального використання

Надалі, доцільно дослідити розвиток LLM у період з 2018 по 2023 р., демонструючи їхню еволюцію як закритих, так і відкритих моделей. На рис. 2.2 наведено детальний опис ключових етапів розвитку цих моделей та їхніх розробників [12].

2018 р. – BERT (Google) – модель, яка революціонізувала підхід до NLP завдяки 2-направленому підходу, що дозволяє краще розуміти контекст. GPT-1 (OpenAI) – перша модель в серії GPT, яка започаткувала розвиток генеративних моделей, здатних створювати зв'язні тексти.

2019 р. – GPT-2 (OpenAI) – значно більша і потужніша версія GPT-1, здатна генерувати зв'язні тексти та вирішувати різні завдання без додаткового навчання. BART (Facebook) – модель, що поєднує автоенкодер і автодекодер

mT5 (Google) – мультимовна версія моделі T5, що підтримує багато мов, що робить її корисною для глобальних задач NLP. GPT-Neo (EleutherAI) – відкрита альтернатива GPT-3, що дозволяє дослідникам використовувати потужні моделі без обмежень. GPT-J (EleutherAI) – подальший розвиток відкритих моделей GPT, забезпечуючи більшу точність і масштабованість.

2022 р. – Chinchilla (DeepMind) – модель, яка оптимізує використання обчислювальних ресурсів, що робить її ефективнішою для практичного застосування. InstructGPT (OpenAI) – версія GPT-3, навчена з урахуванням людських інструкцій, що підвищує її зручність і ефективність у взаємодії з користувачами. WebGPT (OpenAI) – модель, яка використовує інформацію з Інтернету для генерації відповідей, що робить її більш актуальною та інформативною. BLOOM (BigScience) – відкрита багатомовна модель з підтримкою понад 50 мов, що розширює її застосування у глобальному масштабі. PaLM (Google) – модель для глибокого навчання, яка використовує паралельну архітектуру для покращення продуктивності.

2023 р. – GPT-4 (OpenAI) – наступне покоління GPT-моделей з покращеними можливостями, що включає більшу точність і функціональність. Claude (Anthropic) – модель орієнтована на безпечне та етичне використання ІІ, що є важливим для довгострокового розвитку технологій.

LLaMA (Meta) – відкрита модель для різноманітних задач NLP, що робить її доступною для широкого кола дослідників.

В цілому, моделі поділяються на два основні типи:

- відкриті моделі (виділені зеленим кольором): GPT-2 (OpenAI), GPT-Neo (EleutherAI), BLOOM (BigScience), LLaMA (Meta);

- закриті моделі (не виділені кольором): GPT-3 (OpenAI), GPT-4 (OpenAI), Claude (Anthropic), Bard (Google).

Надалі проаналізуємо інформацію про ключових розробників LLM. Компанія Google розробила багато моделей, таких як BERT, T5, PaLM, DeBERTa та ін., що зробили значний внесок у розвиток NLP; к. OpenAI – створила серію GPT-моделей (GPT-1, GPT-2, GPT-3, GPT-4), які значно

вплинули на розвиток генеративних моделей; к. Meta – представила моделі, такі як OPT, LLaMA, які розширили можливості відкритих моделей; к. Baidu – розробила ERNIE, яка підвищила точність розуміння контексту; к. Microsoft – працювала над DeBERTa, Turing-NLG, що покращило якість і точність мовних моделей.

У ШІ GPT-4 від OpenAI є однією з найпотужніших та найвідоміших мовних моделей. Вона здатна виконувати широкий спектр завдань, включаючи генерацію тексту, переклад, відповідь на запитання та багато іншого. Однак, крім GPT-4, існує множина інших альтернатив(рис 2.3), які також пропонують вражаючі можливості і можуть бути використані для різних цілей [15].

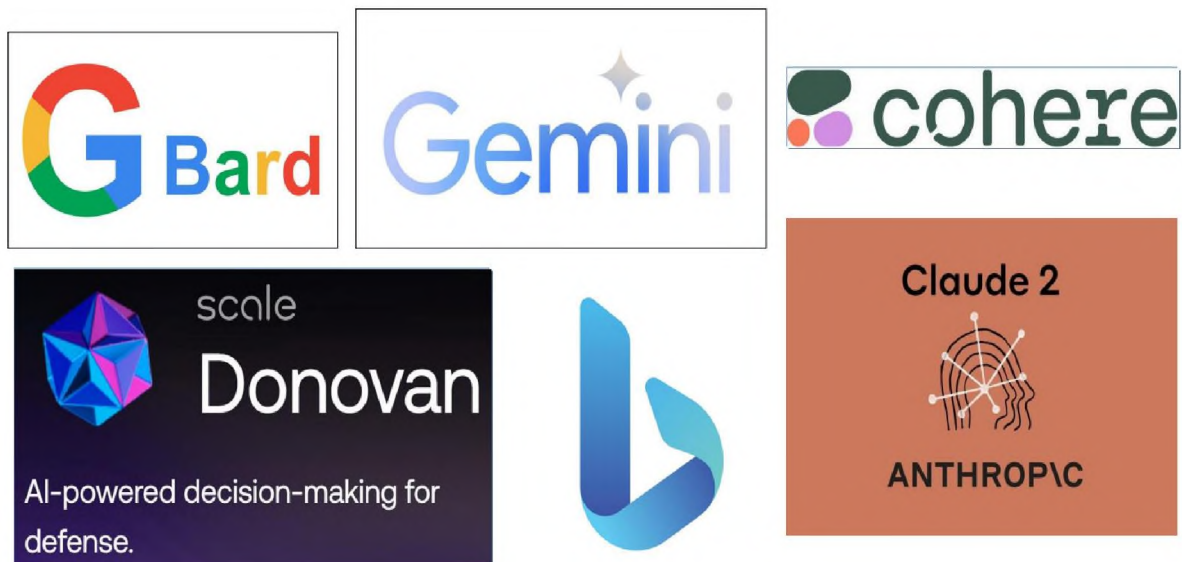


Рисунок 2.3 – Альтернативи GPT-4

В якості прикладу доцільно розглянути Google Bard, також відомий як Gemini. Це чат-бот на основі штучного інтелекту, розроблений компанією Google. Він ґрунтується на LLM LaMDA, а також на новій моделі PaLM 2, яка володіє розширеними можливостями міркування, розуміння мови та кодування. Cohere – LLM розроблена канадською к. Cohere. Вона основана на кількох архітектурах LLM, включаючи Transformer, GPT-3 та Megatron-Turing NLG. Donovan – на жаль, к. Scale AI не розкриває точну інформацію про те, які

LLM використовуються в Donovan. Однак, відомо, що Donovan LLM ґрунтується на кількох різних архітектурах LLM, включаючи Transformer, GPT-3 та Megatron-Turing NLG. Bing, розроблена ко. Microsoft, використовує LLM від OpenAI для покращення своїх функцій NLP та генерації тексту. Зокрема, Bing інтегрує GPT-4 від OpenAI у свій пошуковий механізм та чат, надаючи користувачам більш інтелектуальні відповіді та розширені можливості взаємодії. Claude 3 є самостійною розробкою к. Anthropic (не використовує інші LLM) та створена з урахуванням сучасних досягнень у сфері штучного інтелекту [16].

Найбільшої уваги заслуговують локальні альтернативи GPT-4. Вони дозволяють користувачам запускати потужні інструменти NLP на власних машинах, забезпечуючи контроль над даними, конфіденційність та можливість налаштування під специфічні потреби. Розглянемо основні локальні альтернативи GPT-4, представлені на зображенні (рис. 2.4): Meta LLaMA 2, Stanford Alpaca, Beluga 2, Dolly, Vicuna-13B, GPT4ALL, Baize, Saiga 13B, Falcon 180B, Koala 13B, Mistral 7, 8x7, SynthIA-7B-v1.3, Starling 7B, MPT-30B, Tiny Llama, Gemini та Zephyr [8].

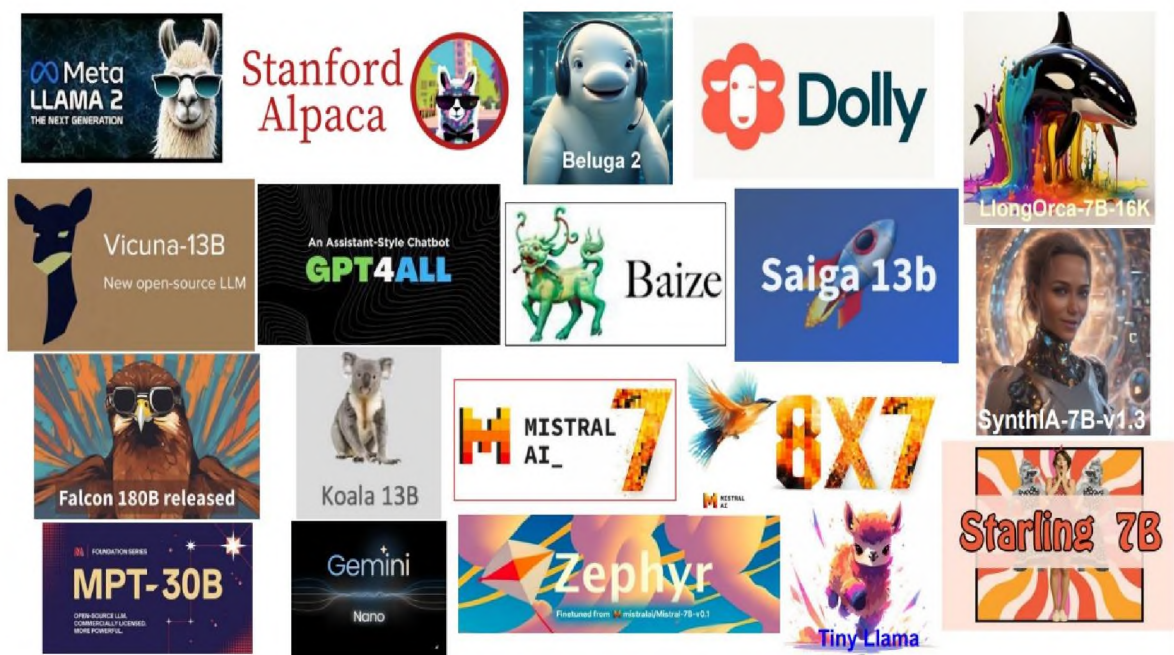


Рисунок 2.4 – Локальні альтернативи GPT-4

В цих моделях є переваги. Контроль над даними – локальні моделі дозволяють забезпечити конфіденційність і безпеку даних, оскільки вони обробляються на власних серверах користувача. Налаштування і оптимізація – можливість налаштовувати моделі під конкретні потреби забезпечує високу ефективність і точність. Автономність – використання моделей локально дозволяє працювати без необхідності постійного підключення до інтернету, що підвищує надійність.

Локальні альтернативи GPT-4 пропонують потужні можливості для обробки природної мови, зберігаючи при цьому гнучкість і контроль, необхідні для багатьох професійних застосувань. Вибір між різними моделями залежить від конкретних вимог користувача та умов застосування, але кожна з представлених моделей може стати ефективним інструментом у різних сферах.

РОЗДІЛ 3

РЕКОМЕНДАЦІЇ ЩОДО РЕАЛІЗАЦІЇ ПЕРСОНАЛЬНОГО АСИСТЕНТА ФЕРМЕРА НА ОСНОВІ ВЕЛИКИХ МОВНИХ МОДЕЛЕЙ

3.1 Формування функціоналу асистента фермера, заснованого на великих мовних моделях

Персональний асистент для фермера, створений на базі LLM, служить потужним засобом для покращення продуктивності аграрних операцій, поліпшення ефективності використання ресурсів і прийняття інформованих рішень. Завдяки інтеграції найсучасніших технологій та ШІ, цей асистент може стати незамінним помічником у роботі кожного дня для фермера [21]. Ось більш детальний опис ключових функцій цього асистента:

1. Аналіз даних. Асистент може обробляти великі обсяги агрономічних даних, включаючи погодні умови, стан ґрунту, рівні вологості та дані про зріст рослин, щоб надавати рекомендації з покращення врожайності та ефективності використання ресурсів.

2. Управління завданнями. Допомогати у плануванні сільськогосподарських операцій, наприклад, у складанні графіка посадки, добрив, поливу та збору врожаю.

3. Діагностика та консультації. Надавати поради щодо лікування рослин, ідентифікувати хвороби чи шкідників на основі фотографій або описів симптомів та рекомендувати методи лікування.

4. Навчання та підтримка. Надавати навчальні матеріали, останні дослідження та кращі практики у галузі сільського господарства, допомагаючи фермерам оновлювати та поглиблювати свої знання.

5. Моніторинг ринку та цін. Аналізувати ринкові тенденції, ціни на сільськогосподарську продукцію та надавати рекомендації щодо оптимального часу для продажу врожаю.

6. Стійкий розвиток. Пропонувати стратегії для стійкого ведення сільського господарства, включаючи методи мінімізації відходів, збереження водних ресурсів та використання альтернативних джерел енергії.

7. Інтеграція з іншими технологіями. Працювати разом з дронами, супутниковими системами та IoT-пристроями для збору даних та автоматизації процесів

8. Мережева взаємодія. З'єднувати фермерів з іншими аграріями, науковими експертами та постачальниками, сприяючи обміну знаннями та ресурсами.

Таким чином, персональний асистент для фермерів на базі LLM – це універсальний інструмент, який надає фермеру якісну підтримку в усіх аспектах сільськогосподарської діяльності, покращує продуктивність ферми і сприяє сталому розвитку аграрного бізнесу.

3.2 Рекомендації щодо реалізації асистента фермера на основі великих мовних моделей

Проектування системи асистента на основі LLM вимагає ретельного планування і реалізації кількох ключових компонентів. Ця система повинна бути здатна збирати та обробляти великі обсяги даних, надавати корисні рекомендації та підтримувати користувача у прийнятті рішень [18]. Нижче наведено короткий опис етапів проектування такої системи.

1. Система асистента повинна мати змогу інтегрувати різні джерела даних, які можуть включати датчики на полі, супутникові зображення, прогностичні дані погоди, інформацію про стан ґрунту та інше.

2. Обробка та аналіз даних, система повинна мати здатність обробляти та аналізувати величезні масиви даних з різноманітних джерел для надання прогнозів про майбутні умови, виявлення проблем та розробки стратегій для їх вирішення.

фермера. Зібрані дані можуть бути використані для подальшого аналізу та прийняття рішень.

Використовуючи LLM від Google Gemini 1.5 Flash, почнемо підключення асистента відносно до Coze.com, а потім до Telegram, таким чином перетворивши його на чат-бота. Для отримання доступу необхідна реєстрація на платформі Coze та отримання API ключів для доступу до сервісів LLM.

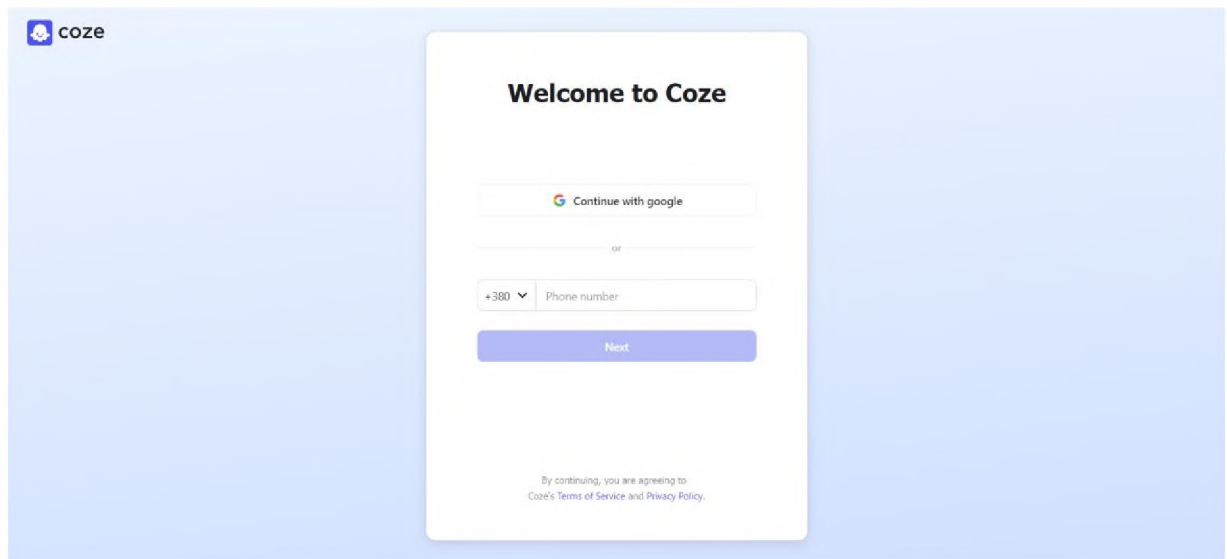


Рисунок 3.2 – Вікно реєстрації на платформі

Coze – це платформа, на якій наразі доступна лише англійська мова. Можливість використовувати українську мову поки що відсутня [19].

Щоб зареєструватися на цьому сайті, ви можете використовувати свій номер телефону або Google акаунт. Процес реєстрації простий і не займає багато часу. Після реєстрації ви отримаєте повний доступ до функцій Coze.

Після успішної реєстрації, ви зможете перейти до створення свого бота. Для цього виберіть «створити нового бота». Тепер потрібно ввести певні деталі про вашого бота. Ці деталі включають ім'я бота, опис, та бажані налаштування (рис 3.3).

Після створення бота, настав час його налаштування. Перше, що потрібно зробити – це вибрати модель LLM, яку потрібно використовувати [31, 32], а також які дії ваш бот виконуватиме (рис 3.4).

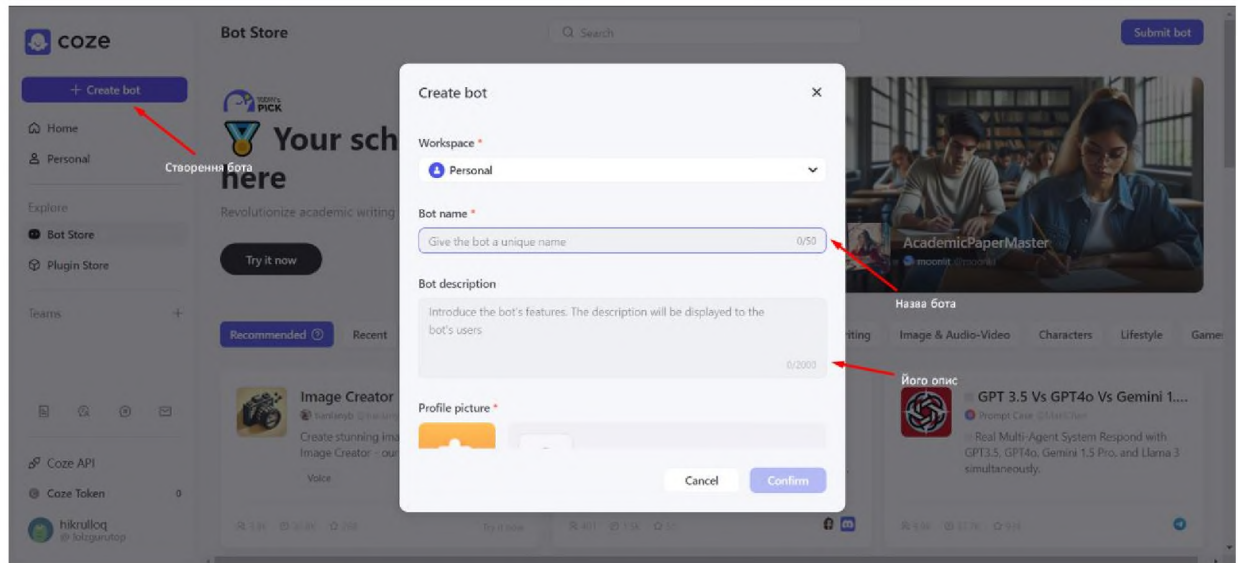


Рисунок 3.3 – Створення бота

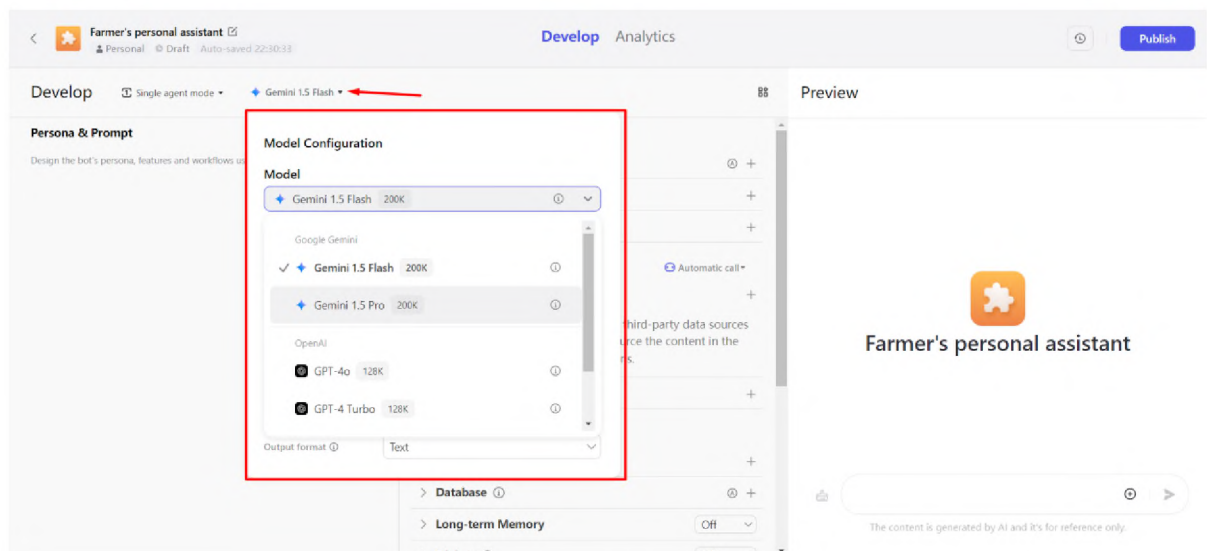


Рисунок 3.4 – Вибір LLM

Платформа Coze наділена величезним набором функцій для налаштування вашого бота, що включає підключення до різних джерел основної інформації. Оскільки опцій багато, давайте сфокусуємось на основних аспектах, які вам необхідно враховувати, для створення персонального асистента фермера на основі різних джерел інформації наприклад книг, веб-сайтів тощо (рис 3.5).

Наповнення інформацією з друкованих джерел:

1. Завантаження вмісту книги, якщо у вас є цифрова версія книги, ви можете використовувати її для наповнення вмістом вашого бота.

2. Ви можете використовувати вбудований інструмент синтаксичного аналізу для розбиття тексту на окремі відповіді та запитання.

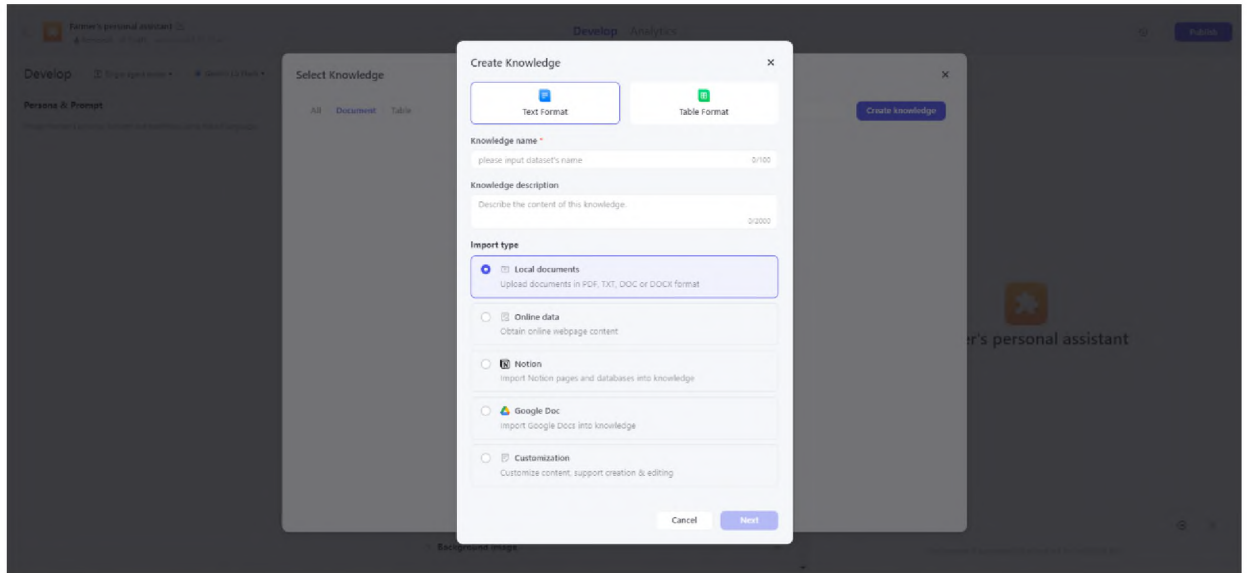


Рисунок 3.5 – Завантаження інформації для обробки

Наповнення інформацією з веб-сайтів:

1. Вебскрейпінг – Coze може запрограмувати вашого бота, щоб він збирав інформацію з декількох веб-сайтів, що транслюють у публічному доступі.

2. Отриману інформацію можна обробити і використовувати в якості відповідей у чаті.

3. Ви можете налаштувати вашого бота таким чином, щоб він автоматично оновлював свою інформацію, збираючи її з вказаних джерел.

Важливо розуміти, що незалежно від обраної вами стратегії, ви повинні завжди пам'ятати про закони, про авторські права і не порушувати їх, надаючи вміст через вашого бота.

Способи включають декілька різноманітних методів імпорту даних у систему, кожен із яких має свої переваги. Перш за все, можна використовувати локальні документи. Це дозволяє завантажувати файли безпосередньо з

комп'ютера користувача у форматах PDF, TXT, DOC або DOCX. Такий підхід ідеальний для роботи з файлами, що вже зберігаються локально, забезпечуючи швидкий та легкий доступ до необхідної інформації.

Інший спосіб – це онлайн дані. Він дозволяє отримувати контент з веб-сторінок. Завантаження тексту або іншої інформації з інтернет-ресурсів для подальшої обробки є корисним для користувачів, яким необхідно працювати з актуальною інформацією з мережі. Також можливий імпорт даних із Notion. Це зручно для тих, хто вже використовує платформу Notion для зберігання своїх даних. Імпорт сторінок і баз даних із Notion дозволяє легко переносити вже збережену інформацію у нову систему.

Інтеграція з Google Документами. Вона дозволяє імпортувати документи із Google Docs, що забезпечує безшовну роботу з хмарними документами, створеними у Google. Налаштування (Customization) дозволяє індивідуалізувати процес завантаження і обробки даних. Це включає підтримку створення та редагування контенту з урахуванням специфічних потреб користувача, можливо, із залученням спеціалістів для адаптації системи.

Ці методи забезпечують користувачам гнучкість і зручність у завантаженні різноманітних типів даних, що дозволяє ефективно використовувати інформацію для подальшої обробки та аналізу.

Після налаштування персонального асистента фермера на платформі Coze, його можна підключити до Telegram для покращення функціональності. Це дозволить отримувати миттєві відповіді на запити та рекомендації в реальному часі, використовуючи зручний інтерфейс месенджера.

Варіант створення бота наведено на рис 3.6.

1. В першу чергу, давайте створимо нового бота в Telegram. Щоб це зробити, потрібно написати BotFather в Telegram, як це робиться зі звичайними контактами, і відкрити чат з ним.

2. Потім потрібно натиснути на команду `"/newbot"` у діалоговому вікні BotFather. BotFather запитає вас про ім'я вашого нового бота і його унікальне ім'я користувача.

3. Після того як ви введете цю інформацію, BotFather надасть вам API токен. Цей токен важливо зберегти, оскільки він потрібен для подальшої інтеграції з Coze.

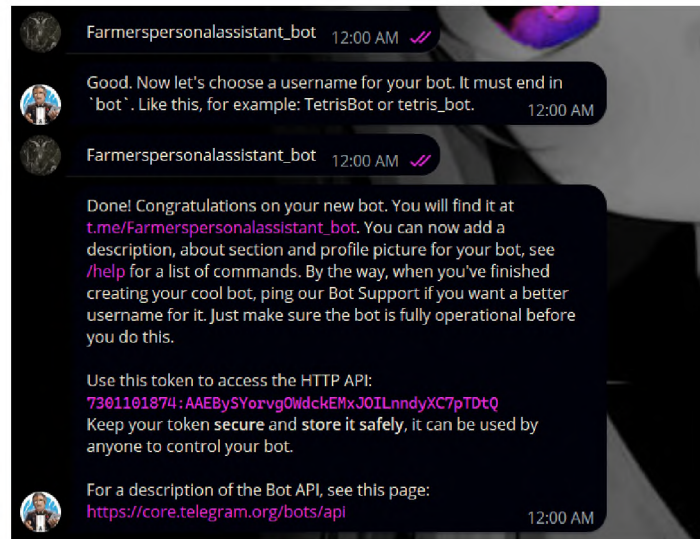


Рисунок 3.6 – Створення бота в Telegram

Далі потрібно пройти процедуру налаштування інтеграції Coze з чат-ботом в Telegram (рис. 3.7).

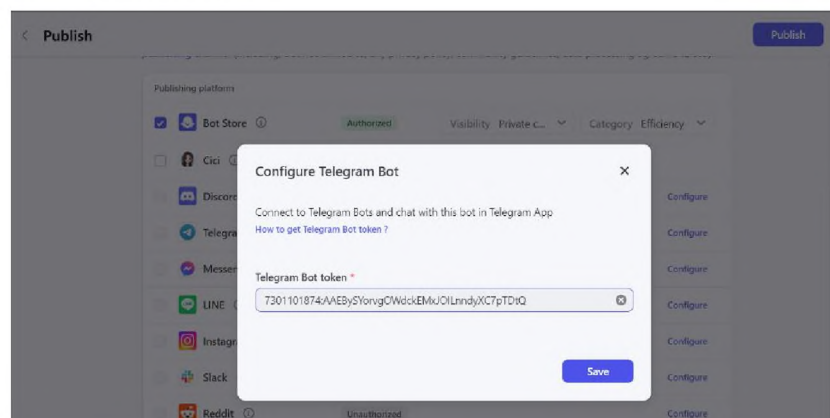


Рисунок 3.7 – Конфігурація API Telegram з Coze

Після інтеграції Coze з Telegram, можна почати його тестувати та використовувати, як персонального асистента фермера.

Тестування бота-помічника для фермера дає можливість оцінити його ефективність в обробці великої кількості інформації і реагуванні на конкретні запити користувача. Бот спрямований на практичну підтримку фермерів, навіть у найвіддаленіших регіонах, за допомогою доступу до сучасних технологій та накопичених знань.

Прогноз погоди є основою для ефективного планування різних аграрних робіт. Бот практично та оперативно надає прогноз погоди для вказаної користувачем локації, допомагаючи планувати роботу на найближчі дні (рис. 3.8).

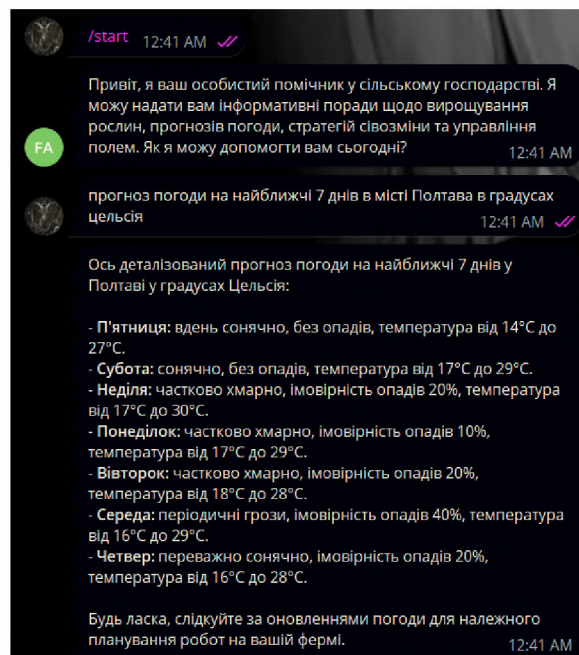


Рисунок 3.8 – Перевірка асистента на прогнозування погодних умов

Бот також забезпечує перевірені та науково обґрунтовані методики сівби культур. Це дозволяє фермерам отримати найновішу інформацію про оптимальні методи вирощування різних культур, включаючи рекомендації щодо густини посадки, обробітку ґрунту, використання добрив та ін. (рис. 3.9).

Складність і виклики, які виникають під час отримання врожаю, можуть бути мінімізовані завдяки боту. Він може надавати конкретні рекомендації для вирішення поточних проблем, а також пропонувати превентивні заходи для запобігання майбутнім проблемам (рис. 3.10).

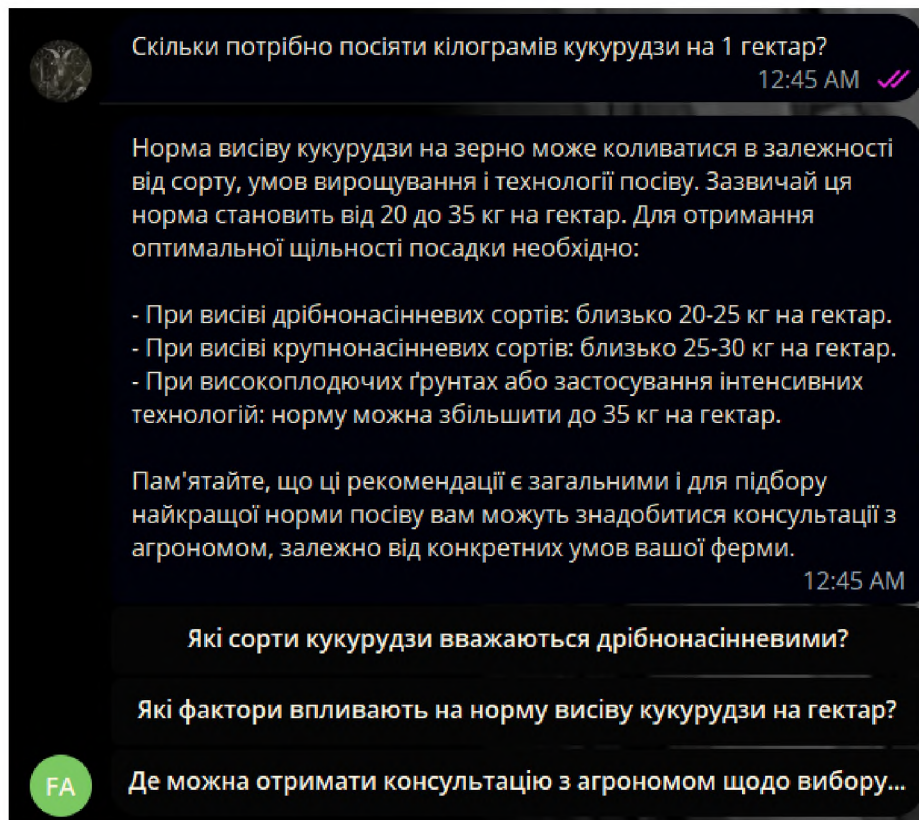


Рисунок 3.9 – Перевірка асистента на методики сівби культур

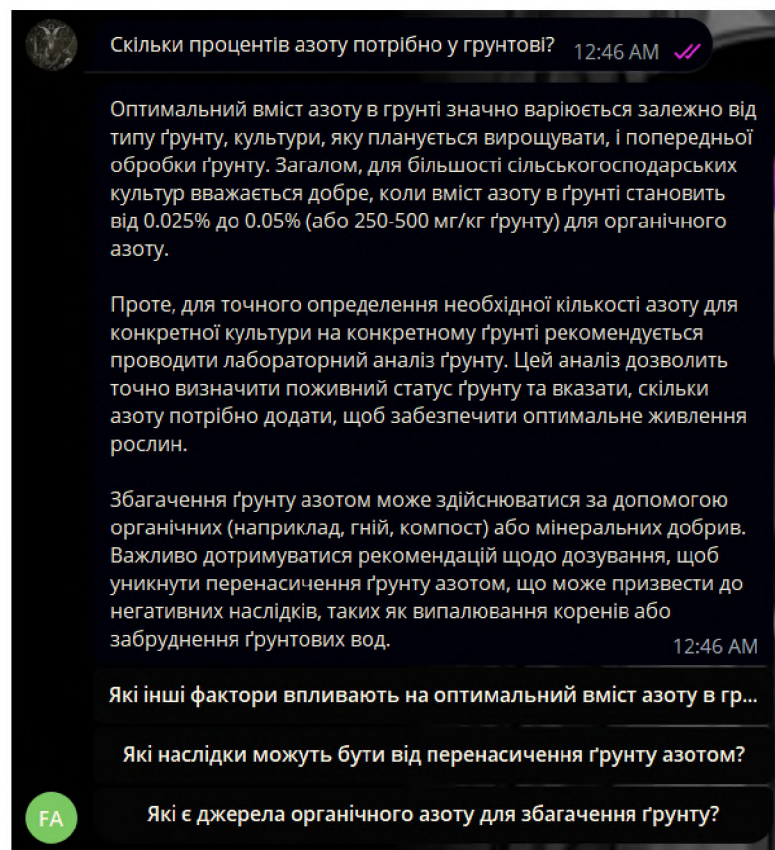


Рисунок 3.10 – Перевірка асистента на рекомендації для вирішення поточних проблем

Все це робить бота цінним інструментом для сучасного фермера, який шукає способи оптимізації своєї роботи та підвищення продуктивності.

3.3 Економічне обґрунтування прийнятих рішень

Економічне обґрунтування є важливим аспектом будь-якого проекту, особливо при інтеграції передових технологій. В даному розділі ми розглянемо витрати та вигоди, пов'язані з розробкою та впровадженням персонального асистента фермера на основі LLM. Проаналізуємо економічну ефективність прийнятих рішень з точки зору інвестиційних витрат, операційних витрат, потенційної економії та зростання продуктивності.

Перш за все, розробка програмного забезпечення вимагає значних інвестицій. Це включає оплату праці розробників, придбання необхідного обладнання. Орієнтовні витрати на початкову розробку включають оплату праці розробників, яка складе приблизно 100000 грн, придбання обладнання – 50000 грн. Таким чином загальні інвестиційні витрати складають 150000 грн. Після впровадження системи необхідно врахувати витрати на її підтримку та обслуговування. Це включає регулярне оновлення програмного забезпечення, забезпечення безперебійної роботи обладнання, навчання персоналу та інші витрати. Річні операційні витрати включають підтримку та обслуговування системи – 20000 грн, оновлення програмного забезпечення – 10000 грн, навчання та підвищення кваліфікації персоналу – 3000 грн, енергетичні витрати – 12000 грн. Загальні операційні витрати складають 45000 грн на рік. Впровадження персонального асистента фермера може призвести до значної економії та вигоди. Основні з них включають підвищення продуктивності праці, зниження витрат на ручну працю, оптимізацію використання ресурсів та підвищення врожайності. Очікується, що підвищення продуктивності праці принесе додаткові 30,000 грн на рік. Зниження витрат на ручну працю складатиме приблизно 15,000 грн на рік.

Оптимізація використання ресурсів зекономить близько 55,000 грн на рік. Підвищення врожайності збільшить дохід на 35,000 грн на рік.

Для оцінки економічної ефективності проекту проводиться аналіз рентабельності. Це включає оцінку терміну окупності інвестицій, розрахунок чистого поточного доходу (NPV) та внутрішньої норми прибутковості (IRR). При цьому NPV розраховується за виразом:

$$\sum_{t=1}^T \frac{C_t}{(1+r)^t} - C_0, \quad (3.1)$$

де C_t – чистий грошовий потік у році t ;

r – дисконтна ставка;

T – загальна кількість років;

C_0 – початкові інвестиційні витрати.

Припустимо, що дисконтна ставка становить 8 %, а очікувані річні грошові потоки від впровадження системи складатимуть 135,000 грн протягом 5 років. За таких умов розрахунок NPV виглядатиме наступним чином:

$$NPV = \frac{135,000}{1.08} + \frac{135,000}{1.1664} + \frac{135,000}{1.2597} + \frac{135,000}{1.3605} + \frac{135,000}{1.4693} - 150,000 = 1188156 \text{ грн.}$$

Таким чином, NPV цього проекту становить 390,348 грн, що свідчить про його економічну вигідність. Оскільки NPV позитивний – проект генерує більше грошових потоків, ніж витрачається на початкові інвестиції та операційні витрати. IRR проекту є значно вищою за дисконтну ставку 8 %, що підтверджує його економічну ефективність. Проведені розрахунки свідчать про значні потенційні вигоди та економію при впровадженні персонального асистента фермера на основі LLM. Хоча початкові інвестиції можуть бути значними, довгострокові вигоди, такі як підвищення продуктивності, зниження витрат та покращення якості продукції, роблять цей проект економічно обґрунтованим та привабливим для інвестицій. Завдяки позитивному NPV та високій IRR, проект є економічно вигідним, що підтверджує доцільність його реалізації.

ВИСНОВКИ

У ході проведених досліджень визначено, що існують різні типи ШІ, такі як машинне навчання, глибоке навчання, обробка природної мови та комп'ютерний зір, кожен з яких має свої унікальні характеристики і сфери застосування. В деяких випадках, ШІ здатний аналізувати дані, вчитися на їхній основі, робити прогнози й ухвалювати рішення з більшою точністю, швидкістю й ефективністю, ніж людина.

ШІ пропонує нові можливості для сільського господарства, допомагаючи вирішувати складні задачі. Наприклад, системи навігації із використанням ШІ можуть допомогти фермерам знаходити найкращі маршрути для поливу та розподілу ресурсів. Інтелектуальні системи можуть керувати високотехнологічною сільськогосподарською технікою, такою як дрони та автоматизовані трактори, що підвищує ефективність роботи та знижує витрати. ШІ також може використовуватись для прогнозування погодних умов, захисту врожаю від шкідників та хвороб, а також для прогнозування попиту на продукти та ринкових цін. Таке застосування ШІ дозволить допомогти фермерам оптимізувати використання ресурсів, знижувати витрати та мінімізувати вплив на навколишнє середовище.

Найбільш перспективною реалізацією ШІ на основі нейронних мереж є архітектура трансформер. Вона базується на концепції само-уваги, є однією з найбільш розповсюджених у завданнях NLP. Трансформер дозволяє обробляти всю вхідну послідовність паралельно, що робить його ефективнішим для довгих послідовностей.

В якості інструментарію для створення персонального асистента фермера розглядаються LLM. На основі оцінки їх властивостей та особливостей застосування розроблено персональний асистент фермера. Він може використовуватися в різних напрямках для підвищення ефективності та спрощення процесів у сільському господарстві. Серед них слід виділити кілька основних:

- аналіз даних;
- управління завданнями;
- діагностика та консультації;
- навчання та підтримка;
- моніторинг ринку та цін;
- стійкий розвиток;
- інтеграція з іншими технологіями;
- мережна взаємодія;

Ці напрямки можуть значно підвищити продуктивність, стійкість та прибутковість сільськогосподарських операцій, допомагаючи фермерам приймати поінформовані рішення та адаптуватися до змінних умов.

Проведене обґрунтування прийнятих рішень за показниками IRR та NPV свідчить про економічну доцільність впровадження персонального асистента фермера (NPV становить 390348 грн). Він може бути реалізований у вигляді мобільного додатку.

Таким чином, результати даної роботи є асистент фермера на основі великих мовних моделей. Він може бути використаний для подальших досліджень за даною тематикою та при проектуванні розумної ферми.