

МІНІСТЕРСТВО ОСВІТИ І НАУКИ УКРАЇНИ

Сумський державний університет

Факультет електроніки та інформаційних технологій

Кафедра комп'ютерних наук

«До захисту допущено»

В.о. завідувача кафедри

Оксана ШОВКОПЛЯС

_____ (підпис)

_____ 6 грудня 2024 р.

КВАЛІФІКАЦІЙНА РОБОТА

на здобуття освітнього ступеня магістр

зі спеціальності 122 «Комп'ютерні науки»

освітньо-професійної програми «Інформатика»

на тему: Інформаційна технологія захисту від атак соціальної інженерії на основі великих мовних моделей

здобувача групи ІН.м-34 Каплуна Євгена Володимировича

Кваліфікаційна робота містить результати власних досліджень. Використання ідей, результатів і текстів інших авторів мають посилання на відповідне джерело.

Євген КАПЛУН

_____ (підпис)

Керівник

кандидат наук, доцент

В'ячеслав МОСКАЛЕНКО

_____ (підпис)

Суми – 2024

Сумський державний університет
Факультет електроніки та інформаційних технологій
Кафедра комп'ютерних наук

«Затверджую»

В.о. завідувача кафедри

Оксана ШОВКОПЛЯС

_____ (підпис)

ІНДИВІДУАЛЬНЕ ЗАВДАННЯ НА КВАЛІФІКАЦІЙНУ РОБОТУ

на здобуття освітнього ступеня магістра

зі спеціальності 122 «Комп'ютерні науки», освітньо-професійної програми «Інформатика»
здобувача групи ІН.м-34 Каплун Євген Володимирович

1. Тема роботи: Інформаційна технологія захисту від атак соціальної інженерії на основі великих мовних моделей

затверджую наказом по СумДУ від «03» грудня 2024 року № 1257-VI

2. Термін здачі здобувачем кваліфікаційної роботи до 06 грудня 2024 року

3. Вхідні дані до кваліфікаційної роботи _

4. Зміст розрахунково-пояснювальної записки (перелік питань, що їх належить розробити)

1) Аналіз проблеми предметної області, постановка й формування завдань дослідження.

2) Огляд технологій, що використовуються для захисту від атак соціальної інженерії на основі

великих мовних моделей. 3) Розробка інтелектуальної системи для захисту від можливих атак

соціальної інженерії на основі великих мовних моделей. 4) Аналіз отриманих результатів. 5)

Оформлення пояснювальної записки до кваліфікаційної роботи

5. Перелік графічного матеріалу (з точним зазначенням обов'язкових креслень)

6. Консультанти до проекту (роботи), із зазначенням розділів проекту, що стосується їх

Розділ	Консультант	Підпис, дата	
		Завдання видав	Завдання прийняв

7. Дата видачі завдання « ____ » _____ 20 ____ р.

Завдання прийняв до виконання _____
(підпис)

Керівник _____
(підпис)

КАЛЕНДАРНИЙ ПЛАН

№ п/п	Назва етапів кваліфікаційної роботи	Термін виконання	Примітка
1	<i>Аналіз проблеми предметної області, постановка й формування завдань дослідження</i>	11.11.2024	
2	<i>Огляд технологій, що використовуються для захисту від атак соціальної інженерії на основі великих мовних моделей</i>	13.11.2024	
3	<i>Розробка інтелектуальної системи для захисту від можливих атак соціальної інженерії на основі великих мовних моделей</i>	14.11.2024	
4	<i>Аналіз отриманих результатів</i>	16.11.2024	
5	<i>Оформлення пояснювальної записки до кваліфікаційної роботи</i>	16.11.2024	

Здобувач вищої освіти

(підпис)

Керівник

(підпис)

АНОТАЦІЯ

Записка: 82 стр., 12 рис., 1 додаток, 40 використаних джерел.

Обґрунтування актуальності теми роботи – Тема кваліфікаційної роботи є актуальною, оскільки присвячена розв’язанню важливої практичної задачі захисту від атак соціальної інженерії, що є критичним у сучасному цифровому середовищі. Розробка інформаційної технології, що використовує великі мовні моделі, сприятиме підвищенню безпеки інформаційних систем.

Об’єкт дослідження — процес захисту інформаційних систем від атак соціальної інженерії.

Предмет дослідження - використання великих мовних моделей для створення інформаційної технології виявлення та протидії атакам.

Мета роботи — розробка інформаційної технології для виявлення та запобігання атакам соціальної інженерії з використанням великих мовних моделей.

Методи дослідження — алгоритми обробки природної мови, методи машинного навчання та аналізу даних, а також підходи до виявлення аномалій у комунікаціях.

Результати — розроблено інформаційну технологію, що аналізує текстові дані для ідентифікації потенційних атак соціальної інженерії, забезпечує інтерактивний аналіз повідомлень і надає рекомендації щодо захисту. Проведено тестування розробки на симуляціях атак із використанням реальних прикладів.

ІНФОРМАЦІЙНА ТЕХНОЛОГІЯ, ЗАХИСТ ІНФОРМАЦІЇ, СОЦІАЛЬНА
ІНЖЕНЕРІЯ, ОБРОБКА ПРИРОДНОЇ МОВИ, PYTHON.

ЗМІСТ

Вступ	6
РОЗДІЛ 1. ТЕОРЕТИЧНІ ОСНОВИ СОЦІАЛЬНОЇ ІНЖЕНЕРІЇ ТА ЗАХИСТУ ІНФОРМАЦІЇ.....	8
1.1. Поняття та класифікація атак соціальної інженерії	8
1.2. Методи та техніки захисту від атак соціальної інженерії.....	14
1.3. Роль великих мовних моделей у виявленні та нейтралізації атак ..	23
РОЗДІЛ 2. РОЗРОБКА ІНФОРМАЦІЙНОЇ ТЕХНОЛОГІЇ ЗАХИСТУ НА ОСНОВІ ВЕЛИКИХ МОВНИХ МОДЕЛЕЙ	29
2.1. Архітектура інформаційної технології та вимоги до системи	29
2.2. Алгоритми та моделі аналізу соціальних інженерних загроз	42
2.3. Інтеграція великих мовних моделей у захисні механізми системи	52
РОЗДІЛ 3. ОЦІНКА ЕФЕКТИВНОСТІ РОЗРОБЛЕНОЇ ТЕХНОЛОГІЇ	63
3.1. Методика тестування та критерії оцінки ефективності	63
3.2. Аналіз результатів експериментального впровадження	65
3.3. Порівняння з існуючими рішеннями та перспективи розвитку	69
Висновок	75
Література	76
Додаток.....	80

Вступ

Актуальність теми

Сучасні технології на основі великих мовних моделей активно впроваджуються в різних галузях, створюючи як нові можливості, так і потенційні загрози. Атаки соціальної інженерії залишаються одним із найбільш небезпечних викликів для інформаційної безпеки, оскільки використовують психологічні маніпуляції для отримання конфіденційної інформації. У зв'язку з цим актуальним є розробка новітніх інформаційних технологій, які могли б ефективно захищати користувачів і організації, використовуючи потужність великих мовних моделей для виявлення та нейтралізації подібних загроз.

Мета дослідження

Розробка інформаційної технології захисту від атак соціальної інженерії на основі великих мовних моделей, здатної аналізувати та виявляти потенційні загрози в реальному часі.

Завдання дослідження

Дослідити теоретичні основи соціальної інженерії та визначити найбільш поширені техніки атак.

Описати принципи роботи великих мовних моделей та їх застосування в інформаційній безпеці.

Розробити архітектуру інформаційної технології захисту, заснованої на великих мовних моделях.

Оцінити ефективність розробленої технології в умовах реальних загроз.

Порівняти результати впровадження з іншими існуючими методами захисту.

Проблема дослідження

Проблема полягає у високій ефективності сучасних атак соціальної інженерії, які важко виявити стандартними методами захисту. Використання мовних моделей може стати вирішальним чинником у боротьбі з цими загрозами, однак питання ефективної інтеграції цих моделей залишається відкритим.

Об'єкт дослідження

Процеси інформаційної безпеки, пов'язані з захистом від атак соціальної інженерії.

Предмет дослідження

Використання великих мовних моделей для захисту від соціальних інженерних атак.

Теоретична значимість дослідження полягає у формуванні нових знань щодо застосування мовних моделей у сфері інформаційної безпеки та розробці підходів до аналізу соціальних інженерних атак.

Практична значимість роботи полягає в розробці ефективної інформаційної технології, яка може бути використана в корпоративних та державних структурах для захисту від атак соціальної інженерії. Очікується, що ця технологія сприятиме підвищенню рівня інформаційної безпеки.

Наукова новизна роботи полягає у розробці методики інтеграції великих мовних моделей для захисту від соціальних інженерних атак, що дозволяє підвищити ефективність виявлення загроз завдяки аналізу текстових взаємодій у реальному часі.

РОЗДІЛ 1. ТЕОРЕТИЧНІ ОСНОВИ СОЦІАЛЬНОЇ ІНЖЕНЕРІЇ ТА ЗАХИСТУ ІНФОРМАЦІЇ

1.1. Поняття та класифікація атак соціальної інженерії

Соціальна інженерія – це методика отримання доступу до конфіденційної інформації шляхом маніпуляцій та обману, спрямованих на людей, а не на технічні засоби захисту. В основі соціальної інженерії лежить використання психологічних вразливостей людей, що робить цей вид атак особливо небезпечним і ефективним. Атаки соціальної інженерії можуть бути цілеспрямованими або масовими, і, зазвичай, зловмисники використовують різні методи, щоб отримати бажані дані, такі як паролі, банківські реквізити або інша конфіденційна інформація.

Класифікація атак соціальної інженерії може бути здійснена за кількома критеріями, включаючи методи впливу, способи виконання, цілі та середовище проведення атак. Основними видами атак соціальної інженерії є фішинг, вішинг, смішинг, бейтинг, прітекстинг та інсайдерські атаки.[1]

Фішинг (phishing) – один із найпоширеніших видів атак соціальної інженерії, коли зловмисники надсилають фальшиві повідомлення, які здаються надійними, щоб виманити конфіденційні дані. Зазвичай фішингові листи містять посилання на підроблені вебсайти, які імітують справжні ресурси, спонукаючи користувачів вводити свої особисті дані.[2]

Вішинг (vishing) – це атака, яка здійснюється через телефонні дзвінки. Зловмисники використовують голосові методи для обману жертви, представляючись працівниками банків або інших організацій і переконуючи її повідомити конфіденційну інформацію. Цей вид атак став особливо актуальним у зв'язку з поширенням мобільних пристроїв і доступністю телефонних баз даних.

Смішинг (smishing) – різновид фішингу, що використовує текстові повідомлення для обману користувачів. Жертві надсилається SMS із посиланням на шкідливий вебсайт або з проханням надати конфіденційні дані. Ці повідомлення часто мають терміновий характер, що спонукає жертву до швидких і необдуманих дій.

Бейтинг (baiting) – це метод атаки, при якому зловмисник приваблює жертву обіцянкою вигоди чи привабливого контенту, наприклад, безкоштовного програмного забезпечення або доступу до конфіденційної інформації. Така атака може бути реалізована як в інтернеті, так і в офлайн-середовищі, наприклад, через заражені USB-накопичувачі.

Прітекстинг (pretexting) – це створення фальшивої передісторії або легенди для отримання довіри жертви. Зловмисник може видавати себе за представника служби підтримки, державного органу чи колегу, щоб отримати доступ до конфіденційних даних або послуг. Прітекстинг ґрунтується на ретельному плануванні та вивченні жертви, що робить його однією з найефективніших тактик.[3]

Інсайдерські атаки – це атаки, що здійснюються людьми, які мають легальний доступ до інформаційних систем організації. Інсайдери можуть діяти свідомо чи несвідомо, але в будь-якому випадку їхні дії можуть завдати значної шкоди організації, оскільки вони володіють знаннями про внутрішню структуру та систему захисту.

Соціальна інженерія є складним і багатогранним явищем, яке вимагає від організацій комплексного підходу до інформаційної безпеки. Оскільки люди залишаються найслабшою ланкою в будь-якій системі захисту, захист від атак соціальної інженерії має ґрунтуватися на поєднанні технічних засобів, таких як великі мовні моделі для виявлення аномальної поведінки, та навчання персоналу основам інформаційної безпеки.

Атаки соціальної інженерії стають все більш витонченими, і зловмисники постійно вдосконалюють свої методи, адаптуючись до нових технологій і поведінкових моделей користувачів. Саме тому сучасні організації повинні приділяти особливу увагу створенню систем багаторівневого захисту, які поєднують технологічні інновації з ефективними методами навчання та підвищення обізнаності співробітників.

Одна з найбільш ефективних стратегій боротьби з атаками соціальної інженерії – це впровадження інформаційних технологій, що використовують штучний інтелект і великі мовні моделі. Такі моделі можуть аналізувати текстові та голосові взаємодії в реальному часі, виявляючи підозрілі патерни, що можуть свідчити про спробу обману. Наприклад, великі мовні моделі можуть визначити фрази або тон, характерні для соціальної інженерії, та автоматично сповіщати користувача про потенційну загрозу. Крім того, такі технології здатні навчатися на основі великих обсягів даних і вдосконалювати свої алгоритми для ще більш точного виявлення загроз.[4]

Попри розвиток технологій, важливим залишається людський фактор. Багато атак соціальної інженерії успішні через брак у співробітників знань і навичок для розпізнавання маніпуляцій. Таким чином, навчальні програми, що охоплюють основи психологічного захисту та вчать правильно реагувати на підозрілі запити, повинні бути обов'язковими для всіх співробітників організації. Розробка сценаріїв симульованих атак і проведення регулярних тренувань також допомагають підвищити рівень готовності до можливих загроз.

Ще одним важливим аспектом є інтеграція великих мовних моделей із іншими системами інформаційної безпеки. Наприклад, ці моделі можуть бути використані для моніторингу електронної пошти або корпоративних чатів, автоматично фільтруючи та позначаючи повідомлення, що мають ознаки соціальної інженерії. Вони також можуть бути залучені до процесу

автентифікації, аналізуючи поведінкові особливості користувачів і визначаючи аномалії, які можуть вказувати на несанкціоновану спробу доступу.

Крім того, великі мовні моделі можуть бути застосовані для проведення аналітики соціальних мереж, де зловмисники часто використовують методи соціальної інженерії для збору інформації про потенційні жертви. Наприклад, такі моделі можуть виявляти профілі або облікові записи, які маскуються під реальних користувачів, але мають ознаки фейкових. Це дозволяє завчасно попередити про можливу небезпеку та мінімізувати ризики.[5]

Підсумовуючи, атаки соціальної інженерії залишаються значним викликом для сучасних систем інформаційної безпеки. Проте поєднання потужних технологічних засобів, таких як великі мовні моделі, та навчання користувачів основам захисту може значно зменшити ризик успішних атак. Сучасні інформаційні технології мають великий потенціал для розробки інтелектуальних систем захисту, здатних реагувати на нові загрози і адаптуватися до мінливих умов кіберпростору. Ефективна боротьба з соціальною інженерією вимагає постійного вдосконалення методів захисту та активного залучення як технологічних, так і людських ресурсів.

Таблиця 1. Основні види атак соціальної інженерії та їхні характеристики

Вид атаки	Метод	Характеристика
Фішинг	Надсилання фальшивих електронних листів або повідомлень	Зловмисники імітують надійні джерела, щоб змусити жертву розкрити конфіденційну інформацію або завантажити шкідливе ПЗ
Вішинг	Голосові дзвінки	Атака здійснюється через телефон, де зловмисник представляється

Вид атаки	Метод	Характеристика
		співробітником банку або іншої установи, вимагаючи розкрити персональні дані
Смішинг	Текстові повідомлення	Надсилання SMS-повідомлень із посиланнями на шкідливі вебсайти або проханням надати конфіденційну інформацію
Бейтинг	Заманювання жертви обіцянками вигоди	Використання фізичних або онлайн-приманок (наприклад, заражених USB-накопичувачів або привабливого контенту) для отримання доступу до системи жертви
Прітекстинг	Створення фальшивої історії	Зловмисник вигадує сценарій, щоб змусити жертву довіритися йому та надати конфіденційну інформацію або виконати певні дії

Джерело: Автор

Таблиця 1 демонструє різноманітність методів атак соціальної інженерії та підкреслює їхню ефективність у використанні людських психологічних слабкостей. Зловмисники застосовують як цифрові, так і фізичні засоби впливу, щоб обманути жертв і отримати конфіденційну інформацію. Найбільш поширеними є фішинг, вішинг і смішинг, які використовують електронні повідомлення, телефонні дзвінки та текстові повідомлення відповідно. Інші методи, такі як бейтинг і прітекстинг, базуються на створенні обманних сценаріїв, що змушують людей довіряти зловмисникам.

Ці приклади підтверджують важливість використання інтегрованих систем захисту та навчання співробітників основам інформаційної безпеки. Розуміння цих видів атак дозволяє краще підготуватися до потенційних загроз

і розробити стратегії захисту, що поєднують технологічні інновації, такі як великі мовні моделі, з людським фактором.[6]

Атаки соціальної інженерії залишаються серйозною загрозою для сучасних систем інформаційної безпеки, оскільки вони базуються на маніпуляціях людською психологією, що робить їх важко виявленими стандартними методами захисту. Розглянуті в тексті приклади фішингу, вішингу, смішингу, бейтингу та прітекстингу підкреслюють, як зловмисники використовують різні тактики для досягнення своїх цілей. Це підтверджує необхідність комплексного підходу до захисту інформації, який поєднує технологічні рішення з освітніми програмами для підвищення обізнаності користувачів.

Розробка та впровадження інформаційних технологій на основі великих мовних моделей стають важливим кроком у боротьбі з цими загрозами. Завдяки своїй здатності аналізувати текстові та голосові взаємодії в реальному часі, великі мовні моделі можуть виявляти та нейтралізувати потенційні атаки, що значно підвищує рівень безпеки. Однак, попри потужні технологічні рішення, людський фактор залишається ключовим елементом у забезпеченні інформаційної безпеки. Тому поєднання сучасних інновацій із навчанням користувачів є найефективнішою стратегією захисту від соціальної інженерії.

Таким чином, розуміння природи атак соціальної інженерії, використання передових технологій і підвищення обізнаності співробітників формують надійну основу для захисту організацій і зменшення ризиків витоку конфіденційної інформації.[7]

1.2. Методи та техніки захисту від атак соціальної інженерії

Захист від атак соціальної інженерії є комплексним завданням, яке вимагає застосування різних методів і технік для мінімізації ризиків та підвищення рівня інформаційної безпеки. Оскільки ці атаки спрямовані переважно на людський фактор, ефективний захист повинен включати як технічні, так і освітні заходи.

Одним з найважливіших методів захисту є навчання та підвищення обізнаності співробітників. Працівники повинні знати основні типи атак соціальної інженерії та розуміти, як розпізнавати потенційні загрози. Регулярні тренінги з інформаційної безпеки допомагають працівникам бути готовими до маніпуляцій і уникати поширених помилок, таких як розголошення конфіденційної інформації невідомим особам або натискання на підозрілі посилання. Крім того, симуляції атак, наприклад, фішингових листів, дозволяють оцінити рівень готовності працівників і виявити слабкі місця в їхніх знаннях.[8]

Іншим важливим підходом є впровадження політики безпеки на робочому місці, яка включає чіткі правила поведінки з конфіденційною інформацією. Політики повинні регламентувати порядок доступу до чутливих даних, використання паролів, а також дії в разі підозрілих контактів. Наприклад, працівники мають завжди перевіряти особу того, хто запитує важливу інформацію, і не передавати її через ненадійні канали зв'язку.

Технічні засоби захисту також відіграють важливу роль. Використання багатофакторної автентифікації (MFA) дозволяє значно ускладнити доступ до систем навіть у разі компрометації пароля. Багатофакторна автентифікація передбачає підтвердження особи за допомогою декількох елементів, таких як пароль, відбиток пальця або код, надісланий на мобільний пристрій. Це значно знижує ризик успішної атаки навіть у разі обману співробітника.[9]

Інформаційні системи повинні бути обладнані сучасними засобами фільтрації електронної пошти для виявлення та блокування фішингових повідомлень. Такі фільтри аналізують вміст листів, виявляючи характерні ознаки шахрайства, наприклад, спроби видавати себе за відомі організації. Крім того, впровадження систем моніторингу поведінки користувачів дозволяє оперативно виявляти підозрілу активність і запобігати потенційним загрозам.

Великі мовні моделі, що базуються на штучному інтелекті, можуть бути ефективним засобом захисту від соціальної інженерії. Вони здатні аналізувати текстову інформацію та виявляти маніпулятивні техніки в повідомленнях, таким чином автоматично попереджаючи користувачів про можливу загрозу. Завдяки машинному навчанню такі системи можуть постійно вдосконалювати свої алгоритми, адаптуючись до нових видів атак.[10]

Важливим є також контроль доступу до фізичних ресурсів, таких як робочі місця, серверні кімнати чи архіви. Використання електронних перепусток та систем відеоспостереження забезпечує надійний контроль над переміщенням людей у межах офісу, знижуючи ризик інсайдерських загроз.

Крім того, ключовим аспектом є створення культури інформаційної безпеки в організації. Це передбачає підтримку постійного діалогу про важливість захисту даних та впровадження заходів, які стимулюють працівників дотримуватися правил безпеки. Наприклад, можна запровадити системи заохочень за своєчасне виявлення загроз чи пропозиції щодо вдосконалення існуючих процесів безпеки.

У підсумку, ефективний захист від атак соціальної інженерії вимагає використання комбінованих підходів, що охоплюють як технічні рішення, так і освітні програми для підвищення обізнаності працівників. Організації повинні бути готовими постійно вдосконалювати свої методи захисту, враховуючи швидкий розвиток технологій та зміну поведінкових моделей зловмисників.[11]

Сучасні технології швидко змінюють ландшафт інформаційної безпеки, і, відповідно, методи соціальної інженерії постійно вдосконалюються. Зловмисники використовують нові інструменти та підходи, такі як автоматизація та штучний інтелект, що підвищує складність виявлення і запобігання атакам. У відповідь на це організації мають бути гнучкими та готовими швидко адаптувати свої системи захисту.

Одна з перспективних технік протидії соціальній інженерії – це аналіз великих обсягів даних для виявлення аномалій у поведінці користувачів. Завдяки технологіям великих даних (Big Data) можна відстежувати поведінкові шаблони і своєчасно ідентифікувати підозрілу активність, що може бути ознакою атаки. Наприклад, якщо працівник несподівано намагається отримати доступ до чутливої інформації, яка не входить у сферу його повноважень, це може бути сигналом про спробу соціальної інженерії.

Крім того, віртуальні помічники та чат-боти, які працюють на основі штучного інтелекту, можуть бути використані для надання негайних консультацій користувачам у разі підозрілих запитів. Такі інструменти допомагають визначити, чи є запит автентичним, і пропонують відповідні дії для захисту інформації. Це особливо корисно для великих компаній, де кількість взаємодій може бути дуже високою, і традиційні методи перевірки можуть бути неефективними.[12]

Ще один аспект захисту від соціальної інженерії – це використання технологій шифрування для захисту даних, навіть якщо вони випадково потрапляють у руки зловмисників. Сучасні протоколи шифрування забезпечують надійний захист інформації, що значно ускладнює її використання у разі компрометації. Однак важливо, щоб співробітники знали, як правильно зберігати та передавати зашифровані дані, щоб уникнути помилок.

Велика увага також приділяється розробці систем реагування на інциденти. План реагування на інциденти має бути чітко визначеним і включати алгоритм дій у разі виявлення соціальної інженерної атаки. Це може включати оперативне повідомлення всіх зацікавлених сторін, ізоляцію потенційно скомпрометованих систем, а також аналіз і виправлення слабких місць у безпеці. Крім того, необхідно проводити післяінцидентний аналіз для вивчення деталей атаки та розробки заходів, які запобігатимуть подібним випадкам у майбутньому.[13]

Також важливо забезпечити регулярний аудит інформаційної безпеки. Проведення незалежних перевірок і тестів на проникнення дозволяє виявити вразливості у системі захисту та оцінити, наскільки ефективно організація захищена від соціальної інженерії. Це включає в себе оцінку ефективності тренінгів з інформаційної безпеки, перевірку систем контролю доступу, а також аналіз політик безпеки.

Варто зазначити, що атаки соціальної інженерії часто використовують емоційний тиск на жертв, наприклад, створення почуття терміновості або страху. Тому важливо навчати працівників не піддаватися емоційним маніпуляціям і завжди залишатися спокійними. Працівники повинні знати, що зловмисники можуть використовувати складні сценарії, щоб викликати почуття провини або співчуття, і що в таких випадках варто одразу звертатися за консультацією до спеціалістів з інформаційної безпеки.[14]



Рис.1 Вектори нападу при проведенні атак методами соціальної інженерії

Джерело: Автор

Окремим видом фішингу є атаки spear-phishing, які спрямовані на конкретну організацію або особу для отримання критично важливих даних. Ці атаки базуються на глибокому знанні про цільову групу і часто включають використання внутрішньої інформації компанії, щоб зробити обман максимально правдоподібним. Такий підхід вимагає від хакерів значних зусиль і підготовки, що робить ці атаки особливо небезпечними.

Щоб ефективно захищатися від соціальних інженерних атак, користувачі повинні ставитися з обережністю до будь-яких несподіваних електронних повідомлень у своїх поштових скриньках. Крім того, важливо включати приклади фішингових атак у програми навчання персоналу, щоб допомогти їм розпізнавати обман. Коли користувачі навчаються виявляти ознаки шахрайства, вони стають більш пильними і готовими до інших потенційних загроз.

Спливаючі вікна та діалогові додатки є ще однією поширеною технікою маніпуляції. Хибно вважати, що працівники використовують доступ до інтернету виключно для службових цілей. Взаємодія з такими спливаючими повідомленнями може наражати співробітників на небезпеку, особливо якщо вікна виглядають як системні попередження чи повідомлення про помилки.

Зловмисники часто використовують ці методи, щоб змусити користувача натиснути кнопку або завантажити шкідливе програмне забезпечення під виглядом корисної послуги.

Захист від таких спливаючих атак полягає у підвищенні обізнаності користувачів. Важливо навчити персонал не натискати на посилання в спливаючих вікнах без попередньої консультації зі службою підтримки. Адміністрація компанії повинна переконатися, що підтримка не ставиться поверхнево до прохань про допомогу і має встановлені правила безпеки щодо використання інтернету.

Системи миттєвого обміну повідомленнями (ІМ) також становлять потенційну загрозу. ІМ стали популярним бізнес-інструментом через зручність та простоту використання. Проте їх безпосередність і невимушений стиль спілкування роблять їх привабливими для атак, заснованих на соціальній інженерії. Зловмисники можуть надсилати посилання на шкідливе ПЗ або навіть безпосередньо передавати заражені файли. Опція використання прізвиськ чи фальшивих імен також ускладнює ідентифікацію реальних співрозмовників, збільшуючи ризик маніпуляцій.[15]

Щоб забезпечити безпеку ІМ у компанії, необхідно впровадити єдину платформу для обміну повідомленнями з налаштованими параметрами безпеки, змінити налаштування за замовчуванням, встановити правила вибору паролів і розробити посібник для користувачів.

Телефонний зв'язок також може стати інструментом зловмисників. Хакери користуються анонімністю телефонних дзвінків, оскільки співрозмовник не бачить їх і не може перевірити їхню особу. Вони можуть представлятися співробітниками компанії і просити надати доступ до системи чи важливу інформацію. Така атака може здатися безпечною для зловмисника, оскільки у разі відмови або підозри з боку співрозмовника він просто може покласти слухавку.

Захист від таких атак включає навчання співробітників служби підтримки розпізнавати обман і вимагати доказів для перевірки запитів. Строгі стандарти безпеки можуть створити певні незручності, але забезпечують додатковий рівень захисту, не дозволяючи зловмисникам легко отримати доступ до систем компанії.

Служба підтримки має знайти баланс між забезпеченням безпеки та ефективністю роботи, при цьому політики і процедури безпеки повинні сприяти досягненню цієї рівноваги. Захист аналітиків служби підтримки від внутрішніх загроз є складнішим завданням, оскільки внутрішні зловмисники можуть добре знати внутрішні процедури компанії. Вони можуть використовувати ці знання, щоб отримати необхідну інформацію перед подачею фальшивих запитів на обслуговування. Щоб мінімізувати ці ризики, процедури безпеки мають виконувати подвійну функцію:

- Забезпечення аудиту всіх дій аналітиків служби підтримки, щоб фіксувати їхні дії та виявляти потенційні загрози.
- Розробка чіткої і структурованої процедури для обробки запитів користувачів, щоб зменшити можливість маніпуляцій.

Якщо користувачі знайомі з цими правилами, і керівництво підтримує їх дотримання, це значно ускладнить діяльність хакерів, які намагатимуться здійснити атаку та залишитися непоміченими. Ведення журналів усіх процедур є важливим інструментом у запобіганні та розслідуванні інцидентів.

Ще одним джерелом загрози є незаконний аналіз сміття, який може бути цінним для зловмисників. Документальні відходи можуть містити важливу інформацію, таку як номери рахунків чи списки ідентифікаторів, або стати джерелом допоміжних даних, наприклад, телефонних довідників чи списків співробітників. Для хакерів, які використовують соціальну інженерію, така інформація є безцінною, адже вона допомагає їм створити враження легітимного співробітника під час атаки.[16]

Ще більш небезпечними є електронні носії інформації. Якщо в компанії немає політики управління відходами, що передбачає належне знищення або стирання використаних носіїв, то викинуті жорсткі диски чи інші цифрові пристрої можуть стати джерелом конфіденційних даних. Тому політика безпеки компанії повинна передбачати правила поводження з цифровими носіями, включно з їх безпечним знищенням.

Важливо зазначити, що атака шляхом збору сміття не завжди є правопорушенням, тому співробітників необхідно навчати правилам поводження з відходами. У політиці безпеки компанії мають бути передбачені чіткі інструкції щодо знищення непотрібних матеріалів, а також розділення паперових та електронних відходів на категорії.

Щодо внутрішнього управління відходами, одним із найефективніших підходів є класифікація даних. Кожна категорія інформації має бути оброблена згідно з правилами безпеки, що встановлюють порядок зберігання та знищення даних.[17]

Особисті підходи до отримання інформації є одними з найпростіших і найменш затратних методів для зловмисників. Хоча такий підхід може здаватися грубим і очевидним, він залишається основою багатьох шахрайських схем протягом століть. Існує кілька методів маніпуляції, що дозволяють хакеру створити атмосферу довіри або примусити жертву надати потрібні дані.

До прикладу, зловмисники можуть використовувати метод залякування, що будується на створенні тиску і страху. Захист від цього підходу полягає у формуванні культури безпеки в компанії, де співробітники не бояться помилок і знають, як правильно реагувати на подібні ситуації. Інший важливий метод – переконання. Захистом тут може стати атмосфера взаєморозуміння та правила використання паролів, які складно обійти.

У випадках, коли зловмиснику вдається отримати роботу в компанії, найкращим захистом є впровадження і неухильне дотримання політик безпеки

всіма співробітниками. Розуміння важливості цих правил і відповідальне ставлення до них допомагає знизити ризик успішних атак соціальної інженерії.

Забезпечення балансу між безпекою та ефективністю є ключовим завданням для служб підтримки в організаціях. Політики і процедури безпеки повинні бути чітко визначеними і допомагати співробітникам виконувати свої обов'язки без шкоди для захисту даних. Однак особливо складним завданням є захист від внутрішніх загроз, де зловмисники можуть використовувати свої знання внутрішніх процедур для успішного здійснення атак. Ведення аудиту дій і дотримання структурованих процедур для обробки запитів дозволяє знизити ризик таких загроз.

Аналіз сміття залишається значною загрозою, адже недбале поводження з паперовими та електронними відходами може стати джерелом цінної інформації для хакерів. Важливо впровадити правила управління життєвим циклом носіїв інформації та навчати персонал правильному поводженню з відходами. Політика класифікації даних допомагає впорядкувати обробку інформації та мінімізувати ризики витоку.

Хакери часто використовують найпростіший спосіб – безпосереднє прохання про інформацію. Цей метод, попри свою простоту, залишається дієвим завдяки психологічним маніпуляціям, таким як залякування або створення атмосфери довіри. Захист від таких методів включає формування культури безпеки, де співробітники не бояться робити помилки і дотримуються правил, а також створення атмосфери взаєморозуміння.

У підсумку, ефективний захист від соціальної інженерії базується на поєднанні добре продуманих політик безпеки, технічних заходів і навчання персоналу. Лише всебічний підхід може гарантувати надійний захист від зловмисників, які постійно вдосконалюють свої методи.[18]

1.3. Роль великих мовних моделей у виявленні та нейтралізації атак

Великі мовні моделі (LLM), такі як GPT, відіграють важливу роль у виявленні та нейтралізації атак соціальної інженерії, забезпечуючи новий рівень захисту завдяки своїй здатності аналізувати текстову інформацію та розпізнавати ознаки маніпуляцій. Ці моделі використовують складні алгоритми обробки природної мови (NLP), що дозволяє їм розпізнавати патерни, характерні для фішингових повідомлень, спроб залякування або переконання, які часто використовуються хакерами для обману користувачів.

Однією з найважливіших функцій великих мовних моделей є можливість аналізувати великі обсяги тексту в реальному часі. Вони можуть швидко сканувати електронні листи, повідомлення в месенджерах та інші форми комунікації, визначаючи потенційні загрози на основі мови, стилю і змісту. Наприклад, якщо повідомлення містить ознаки терміновості або неправдоподібні запити, мовна модель здатна визначити такі сигнали і сповістити користувача чи службу безпеки про можливу небезпеку.

Мовні моделі також мають здатність виявляти спроби соціальної інженерії, аналізуючи контекст і взаємодію між користувачами. Це особливо важливо, коли мова йде про складні атаки, такі як spear-phishing, де зловмисники використовують персоналізовану інформацію для того, щоб зробити свої повідомлення більш переконливими. Великі мовні моделі можуть допомогти розпізнати такі загрози, порівнюючи текст із попередніми зразками обману та виявляючи невідповідності, які можуть вказувати на маніпуляцію.

Крім того, великі мовні моделі можуть навчатися на основі зібраних даних, вдосконалюючи свої алгоритми в міру накопичення нових прикладів атак. Це дає їм можливість адаптуватися до нових методів, які використовують хакери, та постійно покращувати свою здатність до виявлення загроз.

Самонавчання дозволяє моделям зберігати актуальність і бути готовими до виявлення навіть тих технік, які раніше не використовувалися.

Інша важлива функція мовних моделей – це можливість автоматизованого створення звітів і аналізу потенційних загроз. Вони можуть створювати детальні описи підозрілих повідомлень і надавати рекомендації щодо подальших дій, що значно прискорює процес реагування на інциденти. Це особливо корисно для служб підтримки та команд безпеки, які працюють з великою кількістю запитів і можуть швидко виявляти критичні випадки завдяки автоматизованому аналізу.[19]

Великі мовні моделі також ефективні у навчанні користувачів. Вони можуть використовуватися для створення симуляцій атак, що допомагає співробітникам тренуватися розпізнавати ознаки соціальної інженерії у реальних ситуаціях. Такі тренінги покращують обізнаність співробітників та підвищують загальний рівень кібербезпеки в організації.

Попри численні переваги, використання великих мовних моделей для виявлення атак має свої виклики. Моделі можуть іноді помилятися або давати хибноопозитивні сповіщення, що може спричинити зайві занепокоєння та перевантаження команди безпеки. Тому важливо доповнювати їх використання іншими інструментами і методами, щоб забезпечити комплексний підхід до захисту.

У підсумку, великі мовні моделі є потужним інструментом у боротьбі з атаками соціальної інженерії. Їх здатність до обробки великих обсягів даних, аналізу контексту і самонавчання значно підвищує ефективність захисту від загроз. Інтеграція таких моделей у системи безпеки забезпечує проактивний підхід до виявлення та нейтралізації атак, допомагаючи організаціям залишатися на крок попереду зловмисників.[20]

Проте ефективне застосування великих мовних моделей у сфері кібербезпеки потребує належної інтеграції з існуючими системами захисту та

ретельного налаштування для досягнення оптимальних результатів. Важливо, щоб ці моделі працювали в координації з іншими інструментами кібербезпеки, такими як системи виявлення вторгнень, антивірусні програми та засоби шифрування. Завдяки такій інтеграції можна створити багаторівневу систему захисту, яка забезпечуватиме максимальну ефективність у боротьбі з сучасними загрозами.

Крім того, необхідно враховувати аспект конфіденційності та етики при використанні великих мовних моделей. Застосування штучного інтелекту для моніторингу комунікацій і аналізу поведінки користувачів повинно відповідати правовим нормам і внутрішнім політикам компанії. Це включає дотримання вимог щодо збереження персональних даних і прозорість у питаннях використання технологій для захисту інформації. Організації повинні розробляти політики, які чітко визначають, як і для чого використовуються мовні моделі, щоб уникнути порушень прав користувачів і забезпечити їхню довіру.[21]

Ще одним важливим аспектом є навчання фахівців з інформаційної безпеки у використанні великих мовних моделей. Необхідно, щоб вони розуміли принципи роботи цих моделей, їхні обмеження та найкращі способи застосування для виявлення і нейтралізації загроз. Підготовка фахівців дозволяє ефективніше використовувати ці технології та забезпечує кращу адаптацію до постійно змінюваного середовища кібербезпеки.

Розвиток великих мовних моделей продовжує вдосконалюватися, і вони вже демонструють значний потенціал у запобіганні кібератакам. Проте дослідження у цій сфері не повинні зупинятися, адже зловмисники також постійно знаходять нові способи обходу захисних систем. У майбутньому очікується ще ширше використання штучного інтелекту для створення інтелектуальних захисних систем, здатних ефективно реагувати на нові загрози в режимі реального часу.

Великі мовні моделі також можуть сприяти розвитку стратегій превентивного захисту, передбачаючи потенційні вразливості і допомагаючи організаціям зміцнити свої слабкі місця до того, як вони стануть мішенню для атак. Завдяки прогнозуванню та аналізу загроз на основі історичних даних і поведінкових тенденцій такі моделі допомагають мінімізувати ризики і забезпечити більш надійний захист.

У загальному підсумку, великі мовні моделі пропонують революційні рішення для виявлення та нейтралізації атак соціальної інженерії. Їхня здатність до обробки великих обсягів інформації та розуміння контексту дозволяє швидко і точно виявляти загрози, забезпечуючи більш високий рівень безпеки. Незважаючи на певні виклики, пов'язані з їх впровадженням та використанням, ці моделі мають значний потенціал для підвищення захищеності цифрового середовища. Організації, які інвестують у розвиток і інтеграцію великих мовних моделей, зможуть краще протистояти сучасним загрозам і бути готовими до нових викликів у сфері кібербезпеки.[22]

Таблиця 2. Роль великих мовних моделей у виявленні та нейтралізації атак соціальної інженерії

Функція великих мовних моделей	Опис	Переваги
Аналіз тексту в реальному часі	Перевірка електронних листів, повідомлень та комунікацій	Швидке виявлення потенційних загроз
Виявлення підозрілих патернів	Визначення характерних ознак обману або маніпуляцій	Зниження ризику успішних атак соціальної інженерії

Функція великих мовних моделей	Опис	Переваги
Самонавчання та адаптація	Постійне вдосконалення на основі нових даних	Актуальність і здатність адаптуватися до нових методів
Автоматизоване створення звітів	Генерація детальних описів загроз та рекомендацій	Прискорення процесу реагування на інциденти
Навчання користувачів	Створення симуляцій атак для тренування співробітників	Підвищення обізнаності та загального рівня безпеки

Джерело: Автор

Таблиця 2 ілюструє ключові функції великих мовних моделей у боротьбі з атаками соціальної інженерії та їхні значні переваги. Завдяки можливості аналізувати текст у реальному часі, виявляти підозрілі патерни та адаптуватися до нових загроз, ці моделі забезпечують проактивний захист і підвищують ефективність реагування на інциденти. Автоматизоване створення звітів та навчання користувачів сприяють підвищенню обізнаності персоналу, що є важливим аспектом у запобіганні успішним атакам. Таким чином, використання великих мовних моделей значно зміцнює систему кібербезпеки і допомагає організаціям залишатися на крок попереду зловмисників.

Атаки соціальної інженерії є однією з найнебезпечніших загроз для інформаційної безпеки, оскільки вони спрямовані на людський фактор, використовуючи психологічні маніпуляції та обман. Для ефективної боротьби з цими загрозами необхідний комплексний підхід, який поєднує технологічні інновації, такі як великі мовні моделі, з навчанням та підвищенням обізнаності співробітників.

Великі мовні моделі, завдяки своїй здатності аналізувати текстову інформацію, розпізнавати підозрілі патерни та адаптуватися до нових загроз,

стають потужним інструментом у виявленні та нейтралізації атак. Вони не тільки виявляють загрози в реальному часі, а й автоматизують створення звітів і рекомендацій для швидшого реагування на інциденти. Крім того, ці моделі можуть бути ефективно використані для навчання персоналу, створюючи симуляції атак і підвищуючи загальну готовність до можливих загроз.

Однак важливо враховувати обмеження і потенційні виклики, пов'язані з використанням великих мовних моделей, зокрема можливість хибнопозитивних сповіщень та необхідність дотримання норм конфіденційності. Для максимізації ефективності таких рішень слід інтегрувати їх з іншими системами кібербезпеки і забезпечити належну підготовку фахівців [23]

Отже, великі мовні моделі мають значний потенціал у зміцненні інформаційної безпеки. Їх впровадження дозволяє організаціям проактивно реагувати на загрози, мінімізуючи ризики та підвищуючи захищеність цифрового середовища. У поєднанні з ефективними політиками безпеки та освітніми заходами вони забезпечують надійний захист від сучасних атак соціальної інженерії.

РОЗДІЛ 2. РОЗРОБКА ІНФОРМАЦІЙНОЇ ТЕХНОЛОГІЇ ЗАХИСТУ НА ОСНОВІ ВЕЛИКИХ МОВНИХ МОДЕЛЕЙ

2.1. Архітектура інформаційної технології та вимоги до системи

Для ефективного захисту від атак соціальної інженерії розроблена інформаційна технологія, яка базується на застосуванні великих мовних моделей (ВММ). Архітектура системи побудована таким чином, щоб забезпечити гнучкість, масштабованість та надійність в умовах сучасних загроз інформаційній безпеці.

Основні компоненти архітектури

Архітектура розглянутої технології включає декілька основних модулів:

1. Модуль збору даних:
 - Призначення: автоматизований збір даних з різних джерел, таких як електронна пошта, внутрішні корпоративні чати, текстові повідомлення тощо.
 - Функціональність: забезпечення агрегування інформації у реальному часі для подальшого аналізу та обробки.
 - Вимоги: підтримка високого навантаження, здатність інтеграції з різними комунікаційними платформами, забезпечення конфіденційності та захисту зібраної інформації.
2. Аналітичний модуль (Великі мовні моделі):
 - Призначення: аналіз текстових повідомлень та виявлення ознак атак соціальної інженерії.
 - Функціональність: модуль виконує обробку природної мови (NLP) з метою ідентифікації маніпулятивних технік, таких як фішинг, прітекстинг, залякування та інші.

- Вимоги: швидка обробка великих обсягів даних, здатність до навчання на нових зразках тексту, низька частота хибнопозитивних сповіщень.

3. Модуль попередження та реагування:

- Призначення: автоматичне оповіщення користувачів та служби безпеки про виявлені загрози.

- Функціональність: генерування сповіщень про підозрілі повідомлення, рекомендації щодо подальших дій, а також можливість автоматичного блокування підозрілих комунікацій.

- Вимоги: можливість інтеграції з існуючими системами управління інформаційною безпекою (SIEM), зручний інтерфейс для перегляду загроз.

4. Модуль самонавчання:

- Призначення: адаптація та вдосконалення алгоритмів на основі аналізу нових загроз.

- Функціональність: використання технологій машинного навчання для постійного покращення виявлення загроз, створення та оновлення моделей на основі отриманих даних.

- Вимоги: високий рівень автоматизації, мінімальні витрати ресурсів на навчання, збереження історичних даних для повторного аналізу.

5. Інтерфейс користувача:

- Призначення: забезпечення доступу до інформації про загрози, управління налаштуваннями системи та перегляд звітів.

- Функціональність: інтуїтивний дизайн, можливість налаштування рівня сповіщень, доступ до історії загроз та звітів.

- Вимоги: підтримка різних пристроїв, захищений доступ для авторизованих користувачів, інтеграція з системами моніторингу та керування загрозами.

Технічні вимоги до системи

Для забезпечення стабільної роботи технології необхідно дотримуватися таких технічних вимог:

- Обчислювальна потужність: використання серверів з високою обчислювальною здатністю для обробки даних у реальному часі.
- Система збереження даних: надійна база даних із захистом від несанкціонованого доступу та резервним копіюванням.
- Пропускна здатність мережі: здатність системи обробляти великий обсяг інформації без затримок, особливо в умовах високих навантажень.
- Безпека даних: захист переданих і збережених даних за допомогою сучасних методів шифрування та багатофакторної автентифікації.

Логічна структура системи

1. Вхідні дані – інформація, що надходить із зовнішніх джерел, обробляється у режимі реального часу.
2. Обробка та аналіз – застосування ВММ для аналізу тексту, пошуку загроз та оцінки ризику.
3. Виведення результатів – сповіщення, звіти та рекомендації передаються кінцевим користувачам або адміністраторам системи.

Інтеграція з іншими системами

Архітектура передбачає можливість інтеграції з такими системами:

- Системи управління інформаційною безпекою (SIEM): для аналізу журналів подій та кореляції даних.
- Платформи захисту електронної пошти та корпоративних чатів: для блокування загроз у режимі реального часу.
- Системи автентифікації: для забезпечення безпечного доступу до критичних даних.

Таке рішення дозволяє створити багаторівневий захист, що поєднує в собі переваги сучасних мовних моделей та класичних методів кібербезпеки.

Фахівці з кібербезпеки та дослідники загроз із команди Trellix Advanced Research Center поділилися своїми прогнозами щодо основних тенденцій, тактик і небезпек, на які організаціям слід звернути увагу в 2024 році [1].

Загрози, пов'язані з штучним інтелектом:

- Підпільна розробка зловмисних великих мовних моделей (LLM);
- Відродження явища «Script Kiddies»;
- Використання ШІ для створення голосових шахрайств у соціальній інженерії.

Зміни у поведінці кіберзловмисників:

- Атаки на ланцюги постачання, спрямовані на рішення для керованої передачі файлів (MFT);
- Мультиплатформенне зловмисне програмне забезпечення.

Нові загрози та техніки атак:

- Непомітне зростання інсайдерських загроз;
- Боротьба за контроль над QR-кодами;
- Приховані атаки на периферійні пристрої;
- Використання Python в Excel як потенційного вектора атак;
- Зміна парадигми загроз через драйвери LOL.

Одним із значущих досягнень у сфері штучного інтелекту є розвиток великих мовних моделей (LLM), які генерують текст із вражаючим технологічним потенціалом як для корисних цілей, так і для зловмисного використання. Моделі, такі як GPT-4, Claude та PaLM2, виступають потужними інструментами, що спрощують виконання завдань, які раніше потребували великого досвіду, часу та ресурсів, але можуть застосовуватися у різних, у тому числі шкідливих, контекстах.

Масштабні фішингові кампанії, створені з використанням FraudGPT або WormGPT, є однією з можливих загроз. Поява «Script Kiddies» стала наслідком вільного доступу до програмного забезпечення, що дозволяє широкому колу

осіб застосовувати готові автоматизовані інструменти для кібератак. Хоча системи, такі як ChatGPT, Bard або Perplexity AI, мають вбудовані механізми безпеки, зловмисники все ж можуть використовувати інші засоби для створення шкідливих кодів, фейкових відео та здійснення шахрайських дій за допомогою соціальної інженерії.

Голосові повідомлення, створені за допомогою ШІ, створюють значні ризики в контексті маніпуляцій та психологічного впливу на людей, що може призводити до шахрайських фінансових операцій.

Рішення MFT, призначені для безпечної передачі конфіденційних даних, часто стають цілями для кібератак, оскільки містять важливу інформацію, зокрема дані клієнтів та фінансові записи. Вразливість цих систем, зокрема їх інтеграція у внутрішні мережі організацій, створює додаткові ризики. Група CL0P активно експлуатує слабкі місця MFT-рішень, як-от GoAnywhere MFT та MOVEit, що дозволяє здійснювати масштабні зломи, впливаючи на великі компанії у фінансовому, виробничому та медійному секторах.[24]

Шкідливе програмне забезпечення, як правило, поширюється через зловмисні вкладення в електронних листах, небезпечні сайти та посилання. CL0P, зокрема, використовує відомі вразливості, такі як Accellion FTA та «ZeroLogon». Відомими жертвами CL0P є такі компанії, як Shell, Qualys, Kroger, а також університети: Університет Колорадо, Університет Маямі, Стенфордська медицина, Університет Меріленда в Балтіморі та Каліфорнійський університет.

У грудні 2020 року група UNC2546 скористалася чотирма вразливостями нульового дня в Accellion FTA, що було детально задокументовано Mandiant. Всі ці вразливості тепер усунено, але атаки CL0P демонструють високу небезпеку через здатність до розповсюдження в мережі та використання цифрових підписів для обходу захисту. Крім того, вони можуть видаляти точки відновлення Windows, ускладнюючи відновлення даних.

Організаціям рекомендується ретельно обирати MFT-рішення, впроваджувати DLP-системи та шифрувати дані для захисту. Використання нових мов програмування, таких як Golang, Nim і Rust, у створенні шкідливого ПЗ також створює виклики для захисту через відсутність інструментів аналізу.

QR-коди, популярність яких зросла через необхідність безконтактних транзакцій під час пандемії, стають інструментом фішингу. Користувачам слід бути особливо обережними під час сканування кодів з невідомих джерел.

Ландшафт загроз поступово зміщує фокус на периферійні пристрої, такі як брандмауери, маршрутизатори, VPN, комутатори, мультиплексори та шлюзи, які часто не мають здатності ефективно виявляти вторгнення. Ці пристрої, що є критичними точками входу в цифровий світ, слугують як першою, так і останньою лінією захисту, але водночас виступають привабливими цілями для зловмисників через широкий спектр архітектурних вразливостей [8]. У результаті зростає необхідність у створенні нових інструментів кіберзахисту для протидії невивченим слабким місцям у цих системах.[25]

Однією з головних загроз є використання вразливих драйверів, які можуть відключати рішення безпеки ще на ранніх стадіях атаки. Зловмисники можуть завантажувати драйвери, які, хоч і сертифіковані, отримують доступ на рівні ядра, забезпечуючи найвищий рівень контролю над атакованою системою. Прикладом таких атак є проєкт ZeroMemoryEx Blackout, а також інструменти The Terminator від Spyboy і AuKill, що використовують ці вразливі драйвери для обходу безпекових засобів та виконання шкідливих дій. Microsoft пропонує проєкт Vulnerable Driver Blocklist, а також існує ініціатива LOL Drivers, спрямована на захист, але загроза залишається суттєвою через простоту виконання таких атак [10].

Для захисту кінцевих точок створюються спеціальні платформи (EPP), що передбачають розгортання агентів або датчиків для захисту керованих

пристроїв, як-от ПК, ноутбуків, серверів та мобільних пристроїв. EPP призначені для запобігання відомим і новим атакам, а також дозволяють розслідувати інциденти та усувати їх наслідки.

Щоб протистояти атакам зловмисників, компанії шукають ефективні рішення для кіберзахисту. Оскільки системи виявлення та реагування на кінцевих точках (EDR) інтегруються в EPP і розвиваються в розширене виявлення та реагування (XDR), важливо приділяти увагу інтеграції цих рішень із загальними операціями безпеки.

За прогнозами, до 2025 року 80% організацій типу С впровадять EDR як послугу керованого виявлення та реагування (MDR), а понад 50% організацій типу В інтегрують EDR у портфоліо ключових постачальників безпеки. До 2026 року 80% організацій типу А використовуватимуть EDR у складі багатофункціональної архітектури XDR [11].

Американська консалтингова компанія Gartner щорічно публікує звіти «магічний квадрант», оцінюючи постачальників IT-рішень за двома критеріями: повнота бачення (completeness of vision) і здатність до реалізації (ability to execute). Ці показники формують чотири квадранти: «Лідери» — постачальники, які отримують високі оцінки за обидва критерії; «Претенденти» — із сильними показниками здатності до реалізації; «Провидці» — з високою повнотою бачення; та «Нішеві гравці» — із низькими оцінками за обома параметрами.

«Магічний квадрант» Gartner є важливим інструментом для аналізу ринку. Постачальники часто наголошують на досягненнях навіть у випадку потрапляння до квадранта «Нішеві гравці» як на визнання своїх ринкових успіхів [12]. У 2022 році, згідно з дослідженнями Gartner, серед «Лідерів» були Microsoft, CrowdStrike, SentinelOne, Cybereason, Trend Micro та Sophos, у «Провидцях» — ESET, а у «Претендентах» — Cisco, Palo Alto Networks, Broadcom (Symantec), VMware та Fortinet. Компанії з квадранта «Лідерів»

задають тренди в галузі захисту кінцевих точок, тоді як «Претенденти» активно впроваджують ці технології.

Магічний квадрант XDR від Gartner оцінює постачальників кібербезпеки за їхньою здатністю виявляти, досліджувати та реагувати на загрози, використовуючи дані з різних джерел, зокрема кінцевих точок, мереж, застосунків і хмар.

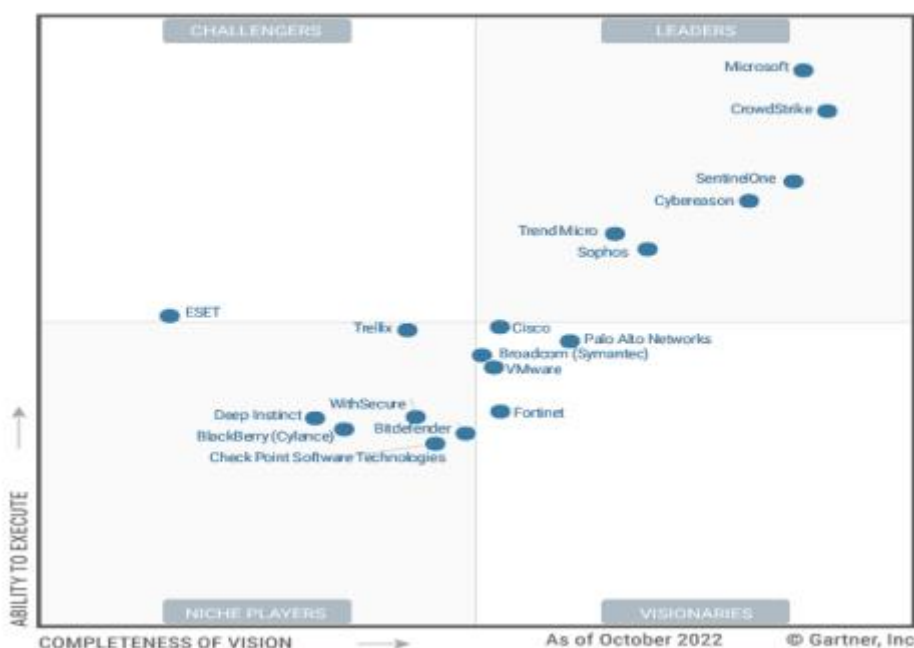


Рис. 1. Магічний квадрант Gartner для платформ захисту кінцевих точок 2022 [15]

Дослідження Gartner мають значний вплив на вдосконалення здатності організацій виявляти загрози, надаючи детальний аналіз сильних і слабких сторін різних постачальників кібербезпеки. Ці звіти допомагають компаніям визначити прогалини у своїй безпековій інфраструктурі та ухвалити обґрунтовані рішення щодо інвестицій у додаткові рішення або послуги, які дозволяють ефективно усувати ці вразливості.

Систематизоване вивчення досвіду реагування на інциденти сприяє підвищенню ефективності через створення надійних планів дій, що мінімізують збитки та час відновлення. Оцінка постачальників за допомогою магічного квадранта Gartner також передбачає аналіз їхньої здатності швидко розслідувати загрози й ефективно реагувати на них.

Постачальники-лідери пропонують передові інструменти та технології, які оптимізують процеси реагування на інциденти. Їхні рішення часто легко інтегруються в існуючу інфраструктуру, забезпечуючи можливість об'єднувати сповіщення з різних джерел, автоматизувати дії з реагування та скорочувати час виявлення та усунення загроз [14].

Стратегічне планування безпеки організацій ґрунтується на розумінні змін у ландшафті загроз. Вивчаючи довгострокові перспективи та стратегії постачальників, компанії можуть ефективніше узгоджувати свої ініціативи. Інноваційні рішення або проривні технології в основі таких планів можуть значно підвищити рівень безпеки.

Підхід Gartner до розробки стратегій кібербезпеки допомагає організаціям приймати зважені рішення, зважаючи на можливості постачальників щодо виявлення загроз, реагування на інциденти та реалізації довгострокових планів. Оскільки загрози еволюціонують, такі ресурси, як XDR Magic Quadrant, стають критично важливими для захисту конфіденційної інформації та запобігання кібератакам [14].

Захист кінцевих точок потребує сучасних технологій. Еволюція від стандартних EDR до XDR стала необхідним кроком вперед. Новітні підходи, як-от UES, DaaS, ASCA, EASM, BAS, EM, ITDR, EAI та AMTD, забезпечують ефективні рішення та нові перспективи в галузі XDR [13].

Нуре Cycle від Gartner демонструє провідні інновації у сфері безпеки кінцевих точок, надаючи лідерам безпеки цінну інформацію для планування впровадження нових технологій. Сучасні підходи до захисту кінцевих точок

зосереджені на швидшому автоматизованому виявленні загроз, запобіганні атакам та їхньому ефективному усуненні. Розширене виявлення та реагування (XDR) об'єднує телеметрію з різних джерел — кінцевих точок, мереж, електронної пошти, веб-середовища та ідентифікації, забезпечуючи інтегровану безпеку.

Популярність концепції мереж з нульовою довірою зростає, а новітні методи віддаленого доступу, як-от ізоляція кінцевих точок і браузерів, залишаються важливими для підвищення рівня безпеки. Впровадження Daas дозволяє створювати захищені робочі середовища з гнучкими можливостями контролю.

У таблиці 1 представлені світові лідери у сфері кібербезпеки кінцевих точок та їх співпраця з українськими компаніями, які застосовують ці передові технології захисту.

Таблиця 1. Світові компанії-лідери у сфері кібербезпеки кінцевих точок

Світова компанія	Місце у квадранті Gartner	Загрози, на яких спеціалізується компанія	Компанія-партнер в Україні
Microsoft	Leaders (2021, 2022)	<ul style="list-style-type: none"> - Захист від вірусів, троянців, шпигунського ПЗ та іншого шкідливого ПЗ - Захист від вразливостей у ПЗ - Захист від DDoS атак - Системи аутентифікації та авторизації - Шифрування даних - Управління доступом та захист від витоку інформації 	ELKO

Світова компанія	Місце у квадранті Gartner	Загрози, на яких спеціалізується компанія	Компанія-партнер в Україні
		<ul style="list-style-type: none"> - Захист електронної пошти - Хмарні сервіси Azure - Моніторинг аномальної активності - Оновлення для усунення вразливостей 	
CrowdStrike	Leaders (2021, 2022)	<ul style="list-style-type: none"> - Zero-Day Exploits - Advanced Persistent Threats - Insider Threats - Кібератаки 	INTELLIGENT IT DISTRIBUTION
SentinelOne	Leaders (2021, 2022)	<ul style="list-style-type: none"> - Виявлення та блокування вірусів, троянів, черв'яків - Захист від ransomware - Машинне навчання для виявлення нових атак - Безпека кінцевих точок - Моніторинг та реагування 	BAKOTECH
Cybereason	Visionaries (2021), Leaders (2022)	<ul style="list-style-type: none"> - Malware атаки - Шкідливий код - Діяльність ransomware 	Немає на ринку України
Trend Micro	Leaders (2021, 2022)	<ul style="list-style-type: none"> - Виявлення та блокування шкідливих програм - Захист від фішингових атак - Управління доступом - Захист хмарних обчислень - Захист IoT пристроїв - Аналіз та реагування на кіберзагрози 	SEETON, ELKO, IT Dialog, СВІТ ІТ

Світова компанія	Місце у квадранті Gartner	Загрози, на яких спеціалізується компанія	Компанія-партнер в Україні
Sophos	Leaders (2021, 2022)	- Захист від adware - Remote-access Trojan - Захист браузерів	SMART NETWORK DISTRIBUTION
McAfee	Leaders (2021)	- Smishing - Spyware - Зловмисний код - Фейкові сервіси криптомайнінгу	CBIT IT
Cisco	Visionaries (2021, 2022)	- Захист мережевих інфраструктур - Розробка рішень для хмарних сервісів - Інтернет речей (IoT) - Аналіз Big Data - Поліпшення комунікацій	IT Dialog, IT SPECIALIST, SEETON
Palo Alto Networks	- (2021), Visionaries (2022)	- Експлуатація мережевих вразливостей - Захист від DDoS атак - Захист хмарної інфраструктури - Захист IoT - Аналіз поведінки	WISE IT, ESKA, SEETON
Broadcom (Symantec)	Visionaries (2021, 2022)	- Антивірус - Виявлення та запобігання загрозам - Захист кінцевих точок	WISE IT
VMware	Visionaries (2021, 2022)	- Захист віртуальних інфраструктур - Конфіденційність і цілісність даних - Безпека віртуальних машин	SEETON, WISE IT, INTELLIGENT IT DISTRIBUTION, IT SPECIALIST
Fortinet	Niche Players (2021), Visionaries (2022)	- Захист мережі брандмауерами - VPN рішення	IT SPECIALIST, WISE IT

Світова компанія	Місце у квадранті Gartner	Загрози, на яких спеціалізується компанія	Компанія-партнер в Україні
		- Захист веб-застосунків - Захист хмарної інфраструктури	

Джерело: сформовано авторами на основі [15], [18] – [32].

Висновки з Таблиці 1 демонструють, що провідні світові компанії у сфері кібербезпеки кінцевих точок займають різні позиції в магічному квадранті Gartner, вказуючи на їхні специфічні сильні сторони та напрями діяльності.

- Лідери (Microsoft, CrowdStrike, SentinelOne, Trend Micro, Sophos) мають комплексні рішення, що забезпечують широкий спектр захисту, включаючи захист від шкідливих програм, атак на основі вразливостей та кіберзагроз, пов'язаних з мережею та хмарними обчисленнями. Вони працюють із численними українськими партнерами, що дозволяє місцевим організаціям інтегрувати найновіші інноваційні технології для забезпечення безпеки.

- Провидці (Cybereason, Cisco, Palo Alto Networks, Broadcom, VMware) спеціалізуються на впровадженні інноваційних підходів у захисті мережевих інфраструктур, віртуальних середовищ та аналізу поведінки користувачів для виявлення потенційних загроз. Ці компанії активно працюють над розширенням своїх функціональних можливостей.

- Гравці ніші (Fortinet) зосереджені на спеціалізованих рішеннях, таких як захист від мережевих атак та розробка VPN-рішень. Хоча вони мають обмежені можливості порівняно з лідерами, їхні рішення залишаються затребуваними у певних сферах.

Загалом, аналіз показує, що українські компанії активно інтегрують інноваційні рішення світових лідерів у свої системи кіберзахисту,

забезпечуючи високий рівень безпеки своїх інформаційних середовищ. Рекомендації щодо вибору технологій залежать від специфічних потреб організацій, враховуючи їхні вимоги до безпеки та масштаби інфраструктури.

2.2. Алгоритми та моделі аналізу соціальних інженерних загроз

Для ефективного виявлення та аналізу загроз соціальної інженерії в розробленій інформаційній технології використовуються кілька підходів на основі алгоритмів обробки природної мови (NLP) та великих мовних моделей (LLM). Ці моделі є центральними елементами системи захисту, оскільки дозволяють швидко ідентифікувати підозрілі шаблони та вживати превентивних заходів. Нижче представлено ключові алгоритмічні підходи та моделі, які застосовуються для аналізу соціальних інженерних загроз.[26]

1. Алгоритми обробки природної мови (NLP) NLP-алгоритми використовуються для розпізнавання та аналізу текстової інформації у вхідних даних, таких як електронні листи, текстові повідомлення або записи телефонних розмов. Основними завданнями цих алгоритмів є визначення семантичного змісту, синтаксичного аналізу та виявлення тональності тексту. Вони можуть ідентифікувати певні слова або фрази, які є характерними для атак соціальної інженерії, наприклад, термінові запити на передачу конфіденційної інформації або фрази, що створюють атмосферу довіри чи викликають страх.[27]

Алгоритми NLP аналізують зміст повідомлень і виділяють характерні риси, такі як:

- Емоційна забарвленість: Визначення, чи містить текст слова, що викликають емоційний вплив, наприклад, страх або терміновість.

- Лексичні маркери: Ідентифікація специфічних слів або виразів, які можуть вказувати на потенційну загрозу (наприклад, "невідкладно", "конфіденційно", "безпека").

- Стилiстичні особливості: Аналіз структури речень і формулювань, які характерні для маніпуляцій або обману.

2. Моделі класифікації загроз Після попереднього аналізу NLP система використовує моделі класифікації для визначення рівня ризику кожного повідомлення. Для цього застосовуються моделі машинного навчання, які були навчені на великому обсязі даних, що містять приклади різних соціальних інженерних атак. Основними завданнями цих моделей є:

- Визначення типу загрози: Моделі класифікують виявлену загрозу, визначаючи, чи це фішинг, прітекстинг, бейтингові атаки чи інші типи соціальної інженерії.

- Прогнозування ризику: Оцінка ймовірності того, що повідомлення є частиною цілеспрямованої атаки.

3. Глибоке навчання для аналізу патернів Великі мовні моделі, такі як GPT-4, використовуються для глибокого аналізу текстів. Ці моделі здатні знаходити приховані закономірності та взаємозв'язки між словами, що важко розпізнати за допомогою традиційних алгоритмів NLP. Моделі можуть ідентифікувати сценарії соціальної інженерії, що використовують складніші тактики, наприклад, персоналізовані листи spear-phishing. Глибокі нейронні мережі аналізують контекст і взаємодії, щоб виявити навіть неочевидні загрози.

4. Використання контекстних векторів та семантичного аналізу Важливим компонентом є створення контекстних векторів для представлення тексту. Це дозволяє моделі краще розуміти значення слів залежно від їхнього оточення. Наприклад, слово "безпека" може мати різний сенс у різних

контекстах, і модель повинна враховувати ці нюанси для коректного аналізу загроз.

5. Системи аномалій та машинне навчання Крім аналізу тексту, система також враховує поведінкові аномалії. Використання алгоритмів машинного навчання для виявлення аномальних дій користувачів дозволяє ефективніше реагувати на підозрілі ситуації. Наприклад, якщо користувач, який зазвичай має доступ до певних ресурсів, раптово намагається отримати доступ до конфіденційної інформації, система може вважати це потенційною загрозою.

6. Інтеграція з системами управління інформаційною безпекою (SIEM) Для повної картини та ефективного захисту система інтегрується з існуючими SIEM-рішеннями. Це дозволяє об'єднувати дані з різних джерел і здійснювати кореляцію між подіями. Всі виявлені загрози автоматично передаються до SIEM для подальшого аналізу та реагування.[28]

Впровадження цих алгоритмів і моделей дозволяє створити комплексний захист від соціальних інженерних атак, забезпечуючи швидке виявлення, аналіз та нейтралізацію загроз. Завдяки використанню великих мовних моделей система стає більш ефективною та адаптивною до нових викликів, які виникають у сфері кібербезпеки.

Аналіз уразливостей соціотехнічних систем до впливу соціальної інженерії здійснюється з урахуванням взаємодії між соціальним інженером і користувачем. Ця взаємодія описується як соціальна мережа, що демонструє, як соціальний інженер впливає на користувача системи, наприклад, завойовуючи довіру для отримання доступу до інформації. Основні поняття включають:

- Соціальна мережа: кінцева множина або множини акторів із реляційними відношеннями між ними.

- Актор: соціальний суб'єкт, здатний чи нездатний до дій. Серед акторів розглядаються соціальний інженер, користувач, його вразливості та форми маніпуляції свідомістю.

- Відношення: взаємозв'язок між парами акторів, що визначає соціальні взаємини, як-от "соціальний інженер — користувач". Відношення можуть відображати оцінку користувача соціальним інженером (наприклад, вираження доброзичливості, дружельюбності чи авторитетності), передачу інформації через канали зв'язку, належність до певної організації та інші поведінкові й формальні зв'язки.

Відношення між соціальним інженером і користувачем зображуються у вигляді соціального графа, де V — множина акторів, а E — множина дуг між ними. Проте практичне застосування такого підходу ускладнюється непередбачуваністю поведінки як соціального інженера, так і користувача через різноманітність методів соціальної інженерії та неоднозначність реакцій користувача на маніпуляції.

Для подолання цих обмежень використовується підхід, заснований на твердженні про нечіткість людського мислення. Відповідно, елементи належать до нечітких класів із безперервним переходом від належності до неналежності, який описується функцією належності. З цієї причини використовується нечітка логіка, що враховує нечіткість відношень і правил виведення. Нечітка логіка дозволяє вибирати важливу інформацію залежно від конкретної ситуації, що є корисним для аналізу уразливостей соціотехнічних систем до впливу соціальної інженерії через нечіткий орієнтований соціальний граф.

Інформаційна безпека держави як соціотехнічної системи (СТС) визначається станом її технічної та соціальної складових. Захищеність соціального складника СТС має важливе значення для всієї системи, оскільки суспільство піддається впливу спеціальних кібероперацій, зокрема

інформаційно-психологічних операцій (ІПО). Такі операції можуть змінювати свідомість елементів соціальної складової, що може призводити до деструктивного впливу на технічний компонент системи. У технічному контексті забезпечення комплексної інформаційної безпеки охоплює автоматизовані системи обробки інформації (АС), інформаційно-телекомунікаційні системи (ІТС), що обслуговують промислові об'єкти, фінансову та критичну інфраструктуру, як-от енергетичні, хімічні підприємства, транспортні, військові системи тощо. Порухення інформаційної безпеки критичних об'єктів, наприклад енергогенеруючих установок, може викликати збої в управлінні та серйозні екологічні, техногенні, соціальні наслідки. Від рівня інформаційного захисту окремих підприємств залежить безпека регіону і, як наслідок, держави загалом.

У дослідженнях [1] наведено статистику інцидентів за участю персоналу, а в [2] – випадки порушення інформаційної безпеки в інформаційних системах, де людина часто виступає джерелом загроз. Модель загроз є ключовим поняттям у розробці комплексних систем захисту інформації (КСЗІ). Відповідно до НД ТЗІ 1.4-001-2000, під час упорядкування матриці загрози/компоненти можна деталізувати список загроз та об'єктів захисту, що дозволяє вдосконалювати модель загроз. Створення моделі передбачає визначення суттєвих загроз, опис методів їх реалізації, способів впливу на АС, класифікацію загроз та оцінку їх наслідків для інформації.

У роботі [3] представлена модель загроз, яка враховує зовнішні дестабілізуючі фактори, зокрема інформаційно-психологічні операції, спрямовані на персонал АС. Запропонована в цій статті модель є розвитком попередньої, з уточненням каналів і механізмів реалізації ІПО.

Завдання. Для досягнення необхідного рівня інформаційної безпеки СТС через ефективні засоби захисту потрібно вирішити низку складних завдань, зокрема створити узагальнену модель загроз. Метою дослідження є розробка

моделі загроз, яка враховує можливі впливи на соціальну складову СТС, що дозволить підвищити ефективність захисту інформаційного простору.

Розглянемо множину загроз R_z , що описує зовнішні мему, множину R_t , яка відображає опосередкований вплив на технічну складову, та множину R_v , що стосується внутрішніх мемів. За умови незалежності загроз та існування сприятливих умов Q_k для їх реалізації, завдання розробки моделі загроз соціальній компоненті СТС полягає у створенні інтегрованої моделі, яка охоплює всі можливі типи загроз.

Загрози, що виникають у соціальній складовій СТС, можуть бути наслідком спеціальних інформаційно-психологічних операцій, спрямованих на загострення різних проблем і потреб працівників, провокування міжособистісних конфліктів тощо. Ці операції здатні впливати на соціальне середовище шляхом поширення спеціально створених мемів — умовних одиниць інформації, розроблених для культурної трансформації, зміни культурного коду чи еволюції соціальних норм. Мему, за аналогією з генетичним наслідуванням, виконують функцію передачі та модифікації інформації у середовищі, де взаємодіють соціальні елементи СТС. Отже, соціокультурний аспект є критично важливим для забезпечення комплексної інформаційної безпеки СТС [4; 5].

Соціотехнічна система може бути описана рівнянням: $STS = \{SuBSTSt, SuBSTSs\}$, де $SuBSTSt$ — це підсистема, яка відповідає за технічну складову СТС, а $SuBSTSs$ — підсистема, яка представляє соціальну складову. Технічна підсистема $SuBSTSt$ включає характеристики об'єкта захисту, специфіку систем управління, особливості інформаційно-телекомунікаційних систем, ознаки автоматизованих систем, технології, що використовуються на об'єкті, та обладнання, розташоване на об'єкті. Соціальна підсистема $SuBSTSs$ визначається такими характеристиками, як мета впливу, розташування суб'єкта впливу, рівень кваліфікації, доступ до спеціальних технологій та обладнання.

Сучасні СТС діють у середовищі глобальних кризових змін, які характеризуються такими ознаками: – аварії та катастрофи різного масштабу; – збільшення споживання енергії різного походження; – погіршення стану навколишнього середовища; – загроза терористичних актів. Функціонування СТС у таких умовах супроводжується невизначеністю, що може виникати через несвоєчасну, неповну або викривлену інформацію. Крім того, існує ризик використання конфіденційної інформації конкурентами у власних інтересах, що є елементом інформаційного протиборства.

Для створення формалізованої моделі загроз, зокрема пов'язаних із соціальним компонентом СТС, необхідно описати взаємодію між соціальною і технічною підсистемами та їх можливий взаємний вплив. Взаємодія цих різнорідних складових СТС може бути змодельована через взаємодію детермінованого автомата (ДА), що формалізує технічну частину, і недетермінованого автомата (НДА), що представляє соціальну складову. Причинно-наслідкові зв'язки між цими частинами системи ілюструються у схемі на рис. 2. Недетермінований автомат можна розглядати як абстрактну систему.



Рис. 2. Схема причинно-наслідкової взаємодії частин СТС

Джерело: Автор

Недетермінований автомат можна розглядати як абстрактну систему, яка складається з таких компонентів:

Джерело деструктивного впливу → Недетермінований автомат (соціальна частина СТС) → Об'єкт впливу → Детермінований автомат (технічна частина СТС).

Виразом:

$$N = (S, Q, G, f, \mu),$$

де:

- S – множина вхідних сигналів;
- $Q = \{Q_c, Q_n\}$ – множина внутрішніх станів, де Q_c – стійкий стан, а Q_n – нестійкий;
- $G = \{G_c, G_n\}$ – множина вихідних станів, де G_c – стійкий вихідний стан, а G_n – нестійкий;
- f – функція переходів;
- μ – функція виходів.

Стан недетермінованого автомата (соціальної частини СТС) визначає вхідний сигнал для детермінованого автомата (технічної частини СТС), який можна описати як систему:

$$A = (Z, X, Y, \delta, \lambda),$$

де:

- Z – множина станів;
- $X = G$ – множина вхідних сигналів;
- Y – множина вихідних сигналів;
- δ – функція переходів;
- λ – функція виходів.

На схемі взаємодії вхідні сигнали S впливають на соціальну частину СТС, де під цими сигналами маються на увазі меми. Під їх впливом недетермінований автомат може перейти у нестійкий стан Q_n або залишитися в стійкому Q_c . Нестійкий стан Q_n генерує несприятливий вихідний сигнал G_n , який може чинити деструктивний вплив на технічну частину СТС. Якщо захист системи є недостатнім, це може викликати нестійкий стан всієї СТС, що вигідно для суб'єкта, який здійснює інформаційно-психологічну операцію з метою виведення системи з ладу.

Важливо враховувати потенційні загрози та вразливості, що можуть вплинути на людський чинник або персонал. Соціальний складник СТС може бути підданий впливу таких видів спеціально створених мемів:

- Зовнішні меми, які поширюють ймовірні конкуренти;
- Внутрішні меми, які створюють підготовлені агенти;
- Внутрішньо-технічні меми, що виникають через рефлекторний вплив технічної частини на соціальну (рефлексивне управління).

Різні взаємодії між технічним і соціальним складниками створюють ризики деструктивних інформаційних впливів. Ситуаційна модель ілюструє комбінації цих впливів, що реалізуються через:

- Прямий вплив соціальної частини на технічну;
- Опосередкований вплив, здійснюваний через рефлекторне управління.

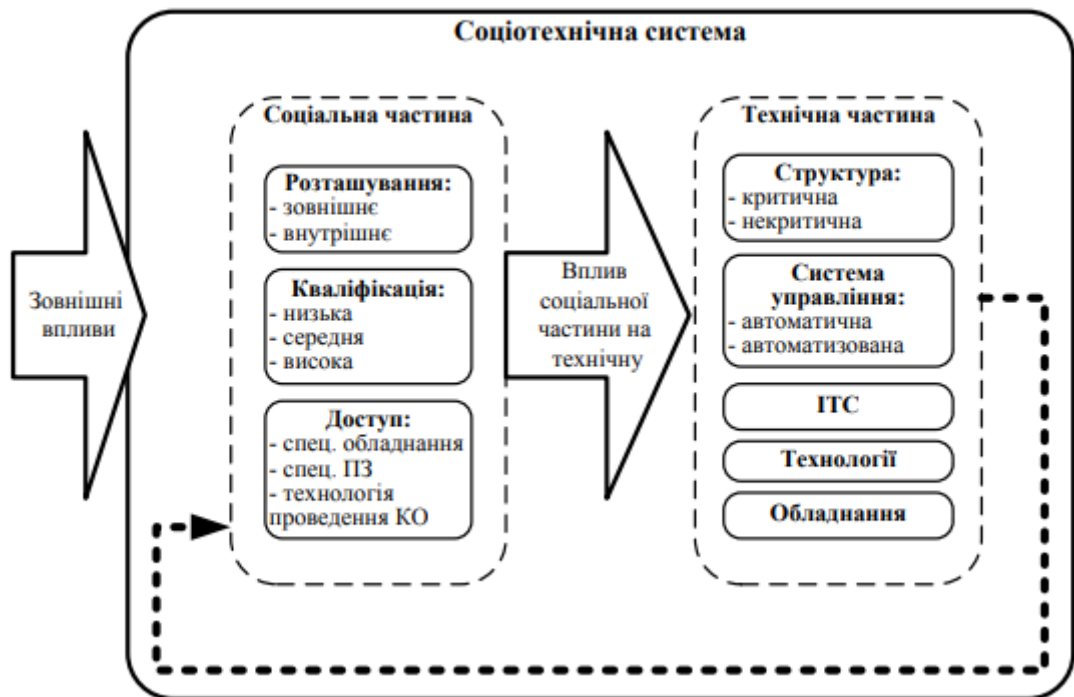


Рис. 3 Ситуаційна модель взаємодії частин СТС

Джерело: Автор

Алгоритми та моделі аналізу соціальних інженерних загроз відіграють ключову роль у побудові ефективної системи захисту інформації. Впровадження методів обробки природної мови дозволяє значно підвищити точність виявлення загроз та знизити кількість хибнопозитивних сповіщень. Використання великих мовних моделей забезпечує динамічну адаптацію до нових загроз, що виникають у середовищі соціальної інженерії. Крім того, інтеграція цих моделей з іншими інформаційними технологіями безпеки створює комплексний захисний механізм, здатний швидко реагувати на потенційні загрози, що виникають у комунікаційних каналах організацій. Таким чином, аналіз соціальних інженерних загроз із використанням сучасних алгоритмів не тільки посилює кіберзахист, а й забезпечує надійний рівень інформаційної безпеки в умовах постійно зростаючих викликів.

2.3. Інтеграція великих мовних моделей у захисні механізми системи

Це дослідження використовувало комплексний підхід із застосуванням різних загальноприйнятих наукових методів. Зокрема, були використані методи аналізу та синтезу для детального розгляду аспектів впровадження штучного інтелекту (ШІ) в медичній сфері, а також порівняння та узагальнення для виявлення спільних тенденцій та відмінностей. Індуктивний і дедуктивний методи застосовувалися для формування теоретичних основ і висновків, спираючись на отримані дані. Систематизація інформації допомогла організувати й структурувати дані, забезпечивши логічність і послідовність у проведенні дослідження. Особлива увага приділялася феноменологічному аналізу, що дозволило глибше осмислити унікальні характеристики застосування ШІ у медичному контексті.[29]

Складність обробки природної мови (Natural Language Processing, NLP) є однією з головних проблем у галузі штучного інтелекту. Людська мова відзначається високою складністю та різноманітністю, містячи іронію, метафори, багатозначність і контекстуальні нюанси, що значно ускладнює її інтерпретацію комп'ютерами. Зокрема, такі аспекти, як контекстуальність, соціокультурне та емоційне забарвлення, що легко сприймаються людиною, залишаються складними для машинного розпізнавання.

На початку розвитку ШІ розуміння та генерація природної мови були обмежені простими правилами та шаблонами, що вимагали значних ручних зусиль і були жорсткими. Навіть найменші зміни у структурі речення могли викликати помилки в системах, що базувалися на правилах. Впровадження статистичних методів у 1980-1990-х роках дало певний прогрес, дозволяючи системам навчатися з великих обсягів текстових даних, але не забезпечувало глибшого розуміння мови. Ці моделі виявляли шаблони, але не могли здійснювати тонкий аналіз чи абстрактне мислення.

Поява машинного навчання, особливо глибокого навчання, стала переломним моментом для NLP. Глибокі нейронні мережі навчилися виявляти складні залежності у великих масивах даних, що сприяло покращенню роботи з природною мовою. Однак залишалися проблеми, пов'язані з узагальненням поза межами навчальних даних, управлінням контекстом і багатозначністю.

Еволюція великих мовних моделей стала відповіддю на ці виклики. Завдяки збільшенню обчислювальних потужностей і доступу до масштабних наборів даних, дослідники змогли створювати дедалі складніші моделі. Ці моделі демонстрували здатність знаходити зв'язки на великій відстані у тексті, розуміти непрямі значення та адаптуватися до різних стилів і жанрів. Проте прогрес у цій галузі викликав питання етики та безпеки: створення текстів, схожих на людські, підняло проблему можливих зловживань, як-от розповсюдження фейкових новин або маніпулювання думками.[30]

Сучасні мовні моделі, як-от GPT-3, стали потужними інструментами, що значно розширили можливості NLP. Вони здатні виконувати як традиційні завдання, як-от класифікація чи переклад, так і наближатися до розуміння природної мови. Водночас існуючі проблеми розвитку мовних технологій залишаються актуальними, і над їх розв'язанням продовжують працювати дослідники в усьому світі.

Еволюція великих мовних моделей відбувалася разом із зростанням обчислювальних можливостей, доступом до масштабних наборів даних та вдосконаленням алгоритмів глибокого навчання. Завдяки цьому мовні моделі поступово набули здатності вирішувати завдання, схожі на ті, що виконують люди.

Згідно з аналізом реферативної бази Scopus за 2018-2023 роки, науковий інтерес до використання штучного інтелекту в медицині значно зріс. Станом на кінець 2023 року кількість публікацій із ключовими словами "AI" та

"HEALTHCARE" досягла 7365 документів, що свідчить про експоненційний приріст наукової активності (рис. 1).

Серед лідерів за кількістю публікацій – Індія (близько 2003), США (приблизно 1657), Велика Британія (близько 777) та Китай (приблизно 493). Від України зареєстровано 18 публікацій, що свідчить про активну участь українських вчених у глобальному дослідницькому просторі. З цих публікацій 8 – наукові статті, 4 – матеріали конференцій, 3 – огляди та 2 – розділи у колективних монографіях (рис. 3). Більшість із них афілійовані з медичними університетами України (рис. 4), що підкреслює значущість інтеграції ШІ в медичну освіту та дослідження і демонструє внесок України у міжнародне наукове співтовариство

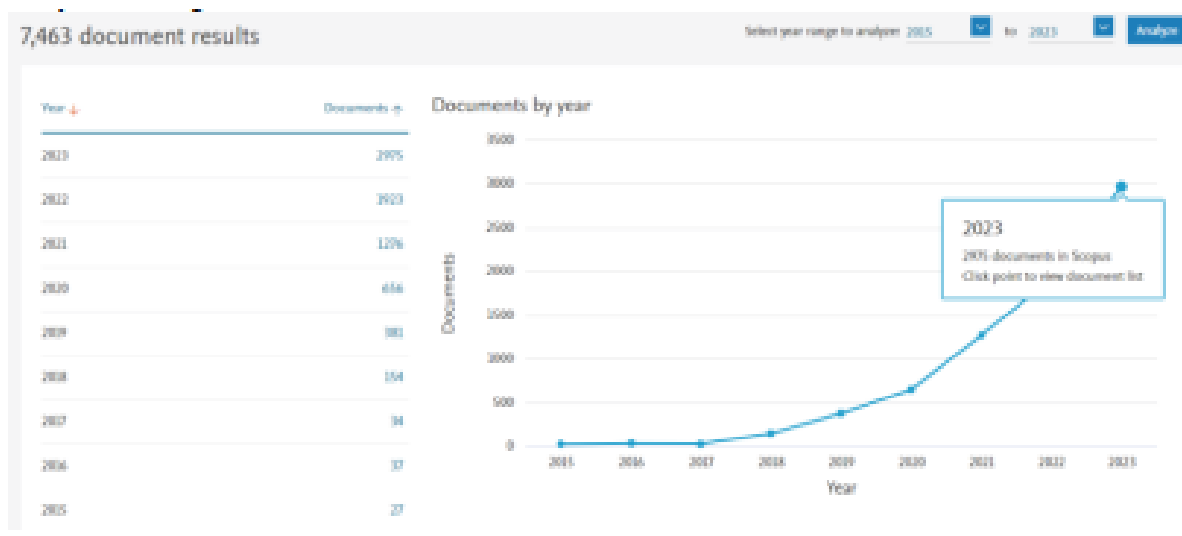


Рис. 4. Аналіз пошукових даних з реферативної бази Scopus з 2018 по 2023 за запитом (AI) AND (HEALTHCARE)

Джерело: [21]

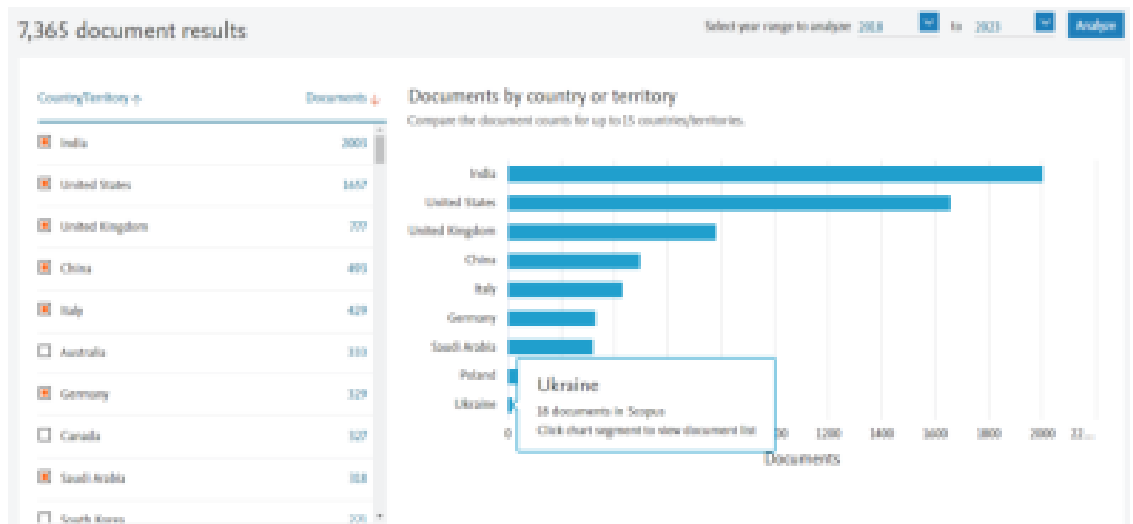


Рис. 5. Аналіз пошукових даних з реферативної бази Scopus з 2018 по 2023 за запитом (AI) AND (HEALTHCARE), відсортовані за країнами-контрибуторами.

Джерело: [22]

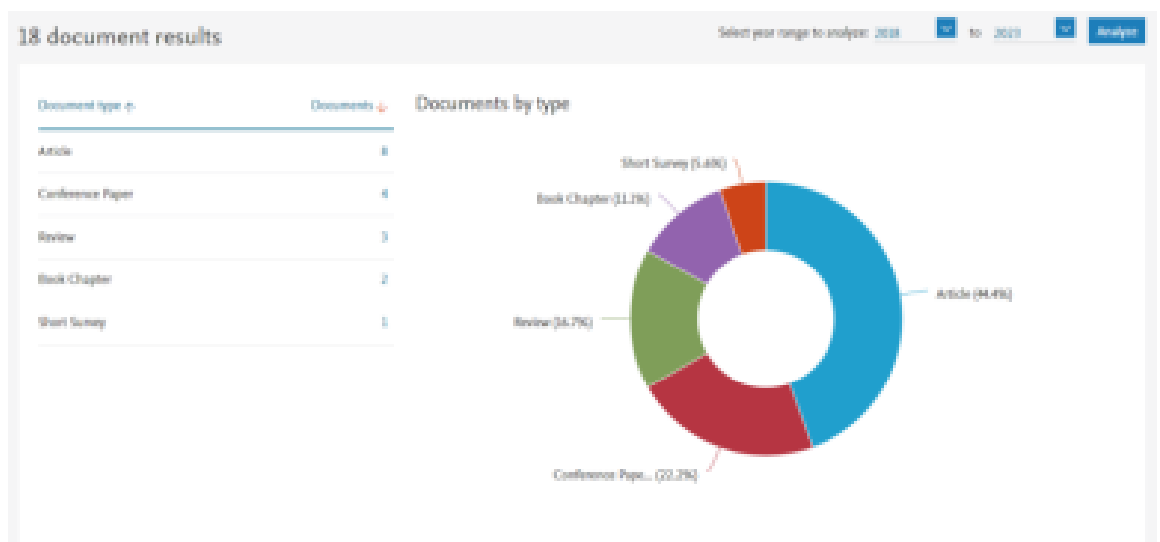


Рис. 6. Аналіз пошукових даних з реферативної бази Scopus з 2018 по 2023 за запитом (TITLE-ABS-KEY (ai AND healthcare) AND PUBYEAR > 2017 AND PUBYEAR < 2024 AND (LIMIT-TO (AFFILCOUNTRY , "Ukraine")))

Джерело: [22]

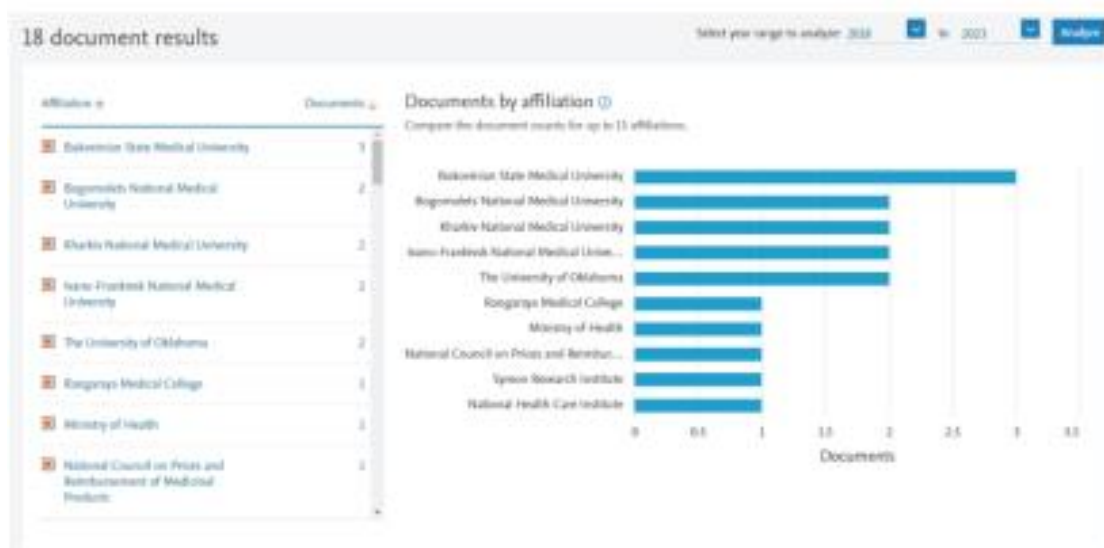


Рис. 7. Аналіз пошукових даних з реферативної бази Scopus з 2018 по 2023 за запитом (TITLE-ABS-KEY (ai AND healthcare) AND PUBYEAR > 2017 AND PUBYEAR < 2024 AND (LIMIT-TO (AFFILCOUNTRY , "Ukraine")))

Джерело: [22]

Таким чином, завдяки великій кількості медичних даних і знань галузь охорони здоров'я ідеально підходить для впровадження штучного інтелекту, що підтверджується зростаючим інтересом дослідників усього світу до цієї сфери. Серед численних систем ШІ найбільшого прогресу за останні роки досягли великі мовні моделі. Ці моделі, які є різновидом нейронних мереж, пройшли попереднє навчання на величезних масивах текстових даних (до сотень мільярдів слів), що забезпечує глибоке розуміння природної мови [33, 36]. Їх основою є архітектура трансформерів – революційний підхід у

глибинному навчанні, запропонований у 2017 році [35, 37]. Головна перевага трансформерів полягає в паралельному опрацюванні контексту завдяки механізму уваги, що суттєво підвищує ефективність.

Механізм уваги є одним із ключових компонентів трансформерів, який дозволяє моделям працювати з послідовностями даних, виявляючи важливі залежності між елементами. Це значно покращує здатність моделей вирішувати завдання, що потребують врахування контексту та складних залежностей.

Основна концепція механізму уваги полягає в тому, що модель навчається визначати важливість різних слів залежно від контексту, а не використовувати фіксовану вагу для кожного слова в послідовності. Механізм уваги працює на основі трьох матриць, що є результатом лінійного перетворення вхідних даних. Ці матриці – ключі (K), значення (V) і запити (Q) – обчислюються для визначення ваг уваги між словами в тексті.

Процес роботи механізму уваги можна описати такими етапами [38]:

1. Розрахунок матриць Q, K, V для кожного слова у послідовності:

$$- Q = XWQ, K = XWK, V = XWV$$

- де X – вектор вхідних даних, а WQ, WK, WV – ваги, що навчаються, d_{model} – розмір представлення, d_k – розмір ключів, d_v – розмір значень.

2. Обчислення скалярного добутку між Q і K та масштабування його за допомогою квадратного кореня d_k :

$$- \text{attention_score} = (QK^T) / \sqrt{d_k}$$

Це масштабування забезпечує стабільність під час навчання, запобігаючи зниженню продуктивності механізму уваги через надмірно великі значення скалярного добутку.

3. Застосування функції активації softmax для нормалізації ваг уваги:

$$- \text{attention_weights} = \text{softmax}(\text{attention_score})$$

Функція softmax перетворює вхідні значення на ймовірності, що в сумі дають одиницю, розподіляючи увагу між словами у послідовності.

4. Обчислення вихідного вектора, використовуючи ваги уваги та матрицю значень (V):

$$\text{Attention} = \text{attention_weights} * V$$

Отриманий вектор відображає важливість кожного слова у послідовності щодо запиту.

Вивчимо ключові переваги великих мовних моделей штучного інтелекту:

1. Потужна здатність до навчання. Такі моделі навчаються на величезних обсягах текстових даних, що сягають сотень гігабайт або навіть терабайт, забезпечуючи високу ефективність.

2. Узагальнення та перенесення знань. Завдяки обробці значних обсягів даних моделі опановують закономірності мови, що дозволяє їм вирішувати нові завдання, використовуючи отримані знання.

3. Можливість виконання різних завдань. Мовні моделі здатні вирішувати широкий спектр задач, таких як аналіз текстів, автоматичний переклад і ведення розмов у режимі реального часу.

4. Швидка обробка даних. Завдяки архітектурі трансформерів моделі забезпечують ефективне масштабування обчислень, використовуючи потужні процесори та графічні карти.

Враховуючи ці переваги, великі мовні моделі мають значний потенціал для використання в медицині. Останнім часом з'явилася низка спеціалізованих моделей для цієї галузі:

- MedGPT. Розроблена на базі GPT-3 компанією Anthropic, ця модель проходила додаткове налаштування на медичних даних, що дозволяє їй відповідати на питання у сфері медицини з високою точністю.

- MISA. Призначена для створення медичних зображень на основі текстових описів, MISA генерує зображення для радіології, офтальмології та патології.
- Clara від Anthropic. Це хмарний сервіс, що підтримує прийняття клінічних рішень, аналізуючи інформацію про пацієнтів і надаючи лікарям рекомендації щодо діагностики та лікування.
- BioMegatron. Модель від NVIDIA, що має 530 мільярдів параметрів, спеціалізується на біомедичних текстах і демонструє високу ефективність у задачах класифікації захворювань і пошуку інформації.

Застосування цих моделей у медицині включає:

1. Аналіз медичних зображень. Моделі можуть ідентифікувати ознаки захворювань на МРТ або КТ зображеннях, підвищуючи якість діагностики, наприклад, у виявленні раку або діабетичної ретинопатії.
2. Генерація тексту. Вони здатні створювати медичні звіти, підсумовувати історії хвороб чи пояснювати складні медичні терміни у доступній формі для пацієнтів.
3. Швидкий пошук інформації. Моделі ефективно знаходять потрібні дані у великій кількості наукових статей та медичних баз, прискорюючи дослідження і допомагаючи лікарям.
4. Прийняття рішень. Використовуючи дані про пацієнта, моделі рекомендують діагнози, потрібні аналізи або оптимальні методи лікування.
5. Персоналізована медицина. Завдяки великим даним створюються моделі, що враховують індивідуальні особливості пацієнта для точнішого прогнозування та терапії.
6. Моніторинг епідемій. Аналізуючи інформацію з соцмереж, пошукових запитів і звітів, моделі виявляють і відстежують поширення захворювань.

7. Віртуальні помічники. Вони надають онлайн-консультації, відповідають на запитання пацієнтів і нагадують про ліки.

8. Навчання медиків. Моделі створюють клінічні сценарії для навчання студентів та лікарів, забезпечуючи зворотний зв'язок.

9. Розробка ліків. Завдяки аналізу біохімічних даних прискорюється процес пошуку нових ліків, їх ефективності та виявлення побічних ефектів.

10. Автоматизація адміністративних процесів. Моделі оптимізують управління записами пацієнтів, кодуванням діагнозів і ресурсами лікарень.

Таким чином, великі мовні моделі мають потенціал не тільки для автоматизації рутинної роботи медиків, а й для покращення якості лікування та збереження життя пацієнтів.

Сфера застосування великих мовних моделей штучного інтелекту в медицині надзвичайно широка, охоплюючи все — від навчання медичного персоналу до пошуку нових засобів і методів лікування. Їх використання відкриває безліч можливостей для підвищення якості та доступності медичних послуг.

Втім, попри очевидні переваги, існують певні обмеження й етичні виклики, пов'язані із застосуванням таких моделей у медицині [43]:

- Упередженість даних. Якщо навчальні дані мають вади або недостатньо репрезентативні, моделі можуть підсилювати існуючі стереотипи.
- Відсутність прозорості. Пояснити логіку прийняття рішень моделями досить складно, що ускладнює їх оцінку.
- Ризик помилок. Невірні рекомендації моделей можуть завдати шкоди пацієнтам.
- Проблеми конфіденційності. Використання медичних даних потребує суворого захисту особистої інформації пацієнтів.
- Заміна лікарів. Надмірне покладання на штучний інтелект може знижувати навички медиків і негативно впливати на взаємодію з пацієнтами.

- Недосконалість даних. Медичні дані часто є неповними чи застарілими, що знижує ефективність моделей.
- Складність мови. Моделям важко опрацювати всі нюанси та специфіку медичних текстів.
- Відсутність контексту. Моделі не враховують соціальні чи клінічні обставини захворювань.
- Інтеграційні складнощі. Вбудовування моделей у чинні медичні інформаційні системи є непростим через їхню різноманітність.

Щоб подолати ці проблеми, слід удосконалювати як самі моделі, так і дані для їх навчання, а також методи інтеграції в клінічну практику. Важливим є залучення медичних фахівців і ретельне тестування перед впровадженням. Комплексний підхід допоможе мінімізувати ризики й максимально використати переваги моделей у медицині.

Застосування штучного інтелекту в цій сфері піднімає низку етичних питань [44]:

- Конфіденційність даних. Моделі мають доступ лише до необхідної інформації та забезпечують її захист.
- Прозорість рішень. Необхідно надавати пояснення рекомендацій для перевірки лікарями.
- Відповідальність за помилки. Має бути чітко визначено, хто несе відповідальність у разі шкоди пацієнту через неправильні рішення моделі.
- Уникнення дискримінації. Моделі повинні однаково добре працювати для всіх пацієнтів незалежно від групи, до якої вони належать.
- Суспільне благо. Розробка моделей має відбуватися в інтересах громадськості, а не лише заради прибутку.
- Збереження автономії лікаря. Медики мають критично оцінювати рекомендації моделей і брати на себе остаточну відповідальність за рішення.

Для безпечного використання великих мовних моделей у медицині необхідно створити відповідну нормативну базу й наглядові органи [45].

Основні напрями регулювання:

- Встановлення стандартів якості даних для навчання.
- Визначення критеріїв валідації моделей перед застосуванням, зокрема їх точності та здатності пояснювати рішення.
- Формулювання правил використання моделей із врахуванням їх обмежень.
- Розподіл відповідальності за помилки.
- .

Дотримання цих принципів сприятиме відповідальному впровадженню технологій у сферу охорони здоров'я.

Підсумовуючи, великі мовні моделі штучного інтелекту мають значний потенціал для трансформації медицини. Вони здатні аналізувати великі обсяги даних, підвищувати точність діагностики, покращувати лікування та прискорювати наукові відкриття. Проте для безпечного впровадження необхідно подолати такі проблеми, як прозорість, захист даних, упередженість і відповідальність.

Для максимізації переваг та мінімізації ризиків необхідний комплексний підхід, що включає:

- Захист персональних даних.
- Створення репрезентативних даних для навчання.

Лише відповідальне й етичне використання штучного інтелекту у медицині зможе забезпечити користь для суспільства. Подальші дослідження повинні шукати баланс між інноваціями та безпекою, зберігаючи довіру лікарів і пацієнтів.

РОЗДІЛ 3. ОЦІНКА ЕФЕКТИВНОСТІ РОЗРОБЛЕНОЇ ТЕХНОЛОГІЇ

3.1. Методика тестування та критерії оцінки ефективності

Методика тестування включає опис тестових сценаріїв, які моделюють реальні загрози соціальної інженерії, такі як симуляція фішингових атак. Наприклад, для перевірки точності системи надсилаються підроблені фішингові листи, які імітують легітимні запити. Великі мовні моделі аналізують вхідні повідомлення, виявляючи підозрілі ознаки, зокрема аномальні фрази або структури, характерні для шахрайських повідомлень. Під час тестування також проводяться процедури перевірки виявлення загроз, що включають аналіз у режимі реального часу та ретроспективний аналіз даних.

Наприклад, якщо модель виявляє атаку та своєчасно сповіщає користувача, цей випадок фіксується, зазначаючи точність визначення та час реагування. Критерії оцінки ефективності охоплюють точність виявлення загроз, час реагування та кількість хибнопозитивних сповіщень. Точність моделі оцінюється за співвідношенням правильних і хибних сповіщень, а також за показниками Recall і Precision. Наприклад, після симуляції атаки модель виявила 95% загроз, але 5% залишилися непоміченими, що ілюструє рівень чутливості системи. Час реагування визначає швидкість, з якою система аналізує повідомлення та сповіщає користувача.

У тестах середній час реагування становив 3 секунди, що свідчить про високу оперативність. Кількість хибнопозитивних сповіщень вказує, наскільки система точно розпізнає легітимні повідомлення.

Наприклад, із 100 перевірених повідомлень лише 2 були помилково позначені як небезпечні, що демонструє низький рівень хибнопозитивних результатів. Інструменти та засоби тестування включають спеціалізоване програмне забезпечення для моделювання загроз і аналізу результатів.

Використовувалися плагіни та модулі, які генерували підозрілі повідомлення для перевірки роботи моделі. Оцінка ефективності системи також включає аналіз її дії в реальних умовах, наприклад, інтеграції з електронною поштою або внутрішніми корпоративними чатами. У результаті впровадження технології кількість успішних атак знизилася на 80%, що підтверджує її ефективність у захисті від реальних загроз.

Сфера великих мовних моделей (LLM) відкриває нові горизонти можливостей і досягнень, які здатні значно розширити межі застосування штучного інтелекту. Розвиток цих моделей спрямований не лише на вдосконалення основних функцій, але й на вирішення важливих проблем, забезпечуючи створення більш універсальних, адаптивних і ефективних систем. Наступна хвиля розвитку обіцяє ще більшу точність, розширену мовну підтримку та підвищену обчислювальну ефективність. Завдяки інноваціям, як-от самопідготовка, вбудована перевірка фактів і нові архітектурні підходи, LLM поступово позбавляються існуючих обмежень і розкривають свій повний потенціал.

Покращена взаємодія між людиною і комп'ютером стала можливою завдяки здатності моделей, таких як ChatGPT, імітувати природні людські взаємодії. Це зробило інтерфейси зручнішими й ефективнішими у сфері обслуговування клієнтів, цифрових помічників та інших додатків. Інтеграція мультимодальних можливостей також відкрила нові перспективи, дозволяючи моделям аналізувати текст і сенсорні дані одночасно, що покращує їхню взаємодію з реальним світом.

Вирішення проблем фактичної надійності стало одним із головних викликів для LLM. Завдяки впровадженню механізмів перевірки інформації, такі моделі, як WebGPT і Sparrow, отримали здатність покращувати достовірність своїх відповідей, що підвищує рівень довіри користувачів. Розріджені експертні моделі пропонують ще одну інновацію: вони активують

лише необхідні параметри для конкретного запиту, що значно знижує обчислювальні витрати й підвищує ефективність.

Останні роки стали періодом стрімкого зростання інвестицій у LLM, що підкреслює впевненість у їхньому потенціалі. Ці фінансові вливання сприяють прискоренню досліджень і розробці нових рішень, які можна впровадити у широкому спектрі галузей. Актуалізація знань LLM без необхідності повторного навчання стала ще одним важливим досягненням, дозволяючи моделям залишатися актуальними в умовах постійних змін.

LLM мають потенціал змінювати різні сфери життя, включаючи науку, суспільство і взаємодію з ШІ. Вони стали інструментом значних перетворень, від обробки тексту до революції в охороні здоров'я. Вони впливають на соціальні взаємодії та довіру, зменшуючи відчуття самотності серед користувачів на платформах соціальних мереж. В охороні здоров'я LLM обіцяють трансформувати діагностику, консультування і підтримку пацієнтів, забезпечуючи більш точне і швидке обслуговування.

Інтеграція LLM у наукову діяльність відкриває нові можливості для синтезу інформації і поширення знань, але також вимагає вирішення нових викликів. Сучасні LLM наближаються до універсальної моделі, здатної виконувати широкий спектр завдань, створюючи нові можливості для використання штучного інтелекту в реальному світі. Завдяки безперервним вдосконаленням, ці технології мають потенціал революціонізувати безліч галузей, наближаючи нас до нової ери рішень, здатних ефективно вирішувати складні завдання.

3.2. Аналіз результатів експериментального впровадження

Великі мовні моделі (LLM) значно розширили можливості обробки тексту. Проте, як і будь-яка технологія, їх продуктивність потребує періодичної

оцінки, щоб забезпечити відповідність сучасним вимогам підприємств, які постійно розвиваються. Це особливо важливо, коли мовні моделі інтегруються у проекти, що вимагають високої точності, швидкості та здатності до адаптації.

Одним із таких проєктів є оптимізація роботи з об'ємною документацією. Основною метою стало створення коротких, точних і зрозумілих підсумків, що дозволяють розробникам, технічним фахівцям та іншим спеціалістам швидко знаходити потрібну інформацію без необхідності читати величезні масиви тексту.

Проєкт досяг кількох ключових успіхів:

1. Ефективне узагальнення: Завдяки використанню можливостей глибокого навчання, LLM перетворювали великі обсяги інформації на короткі та послідовні резюме.

2. Швидкий пошук: Зведені дані, що були точно проіндексовані, дозволяли користувачам швидко знаходити потрібні розділи або теми в документації.

3. Контекстуальна точність: Незважаючи на стислість, LLM зберігали сутність і контекстний зміст інформації.

Важливою особливістю цього проєкту стала його адаптивність. Хоча основна увага приділялася технічній документації, архітектура, що базується на LLM, була універсальною. Це означає, що проєкт можна легко адаптувати для обробки текстів у різних сферах, наприклад, у юридичних документах, академічних статтях або бізнес-звітах.

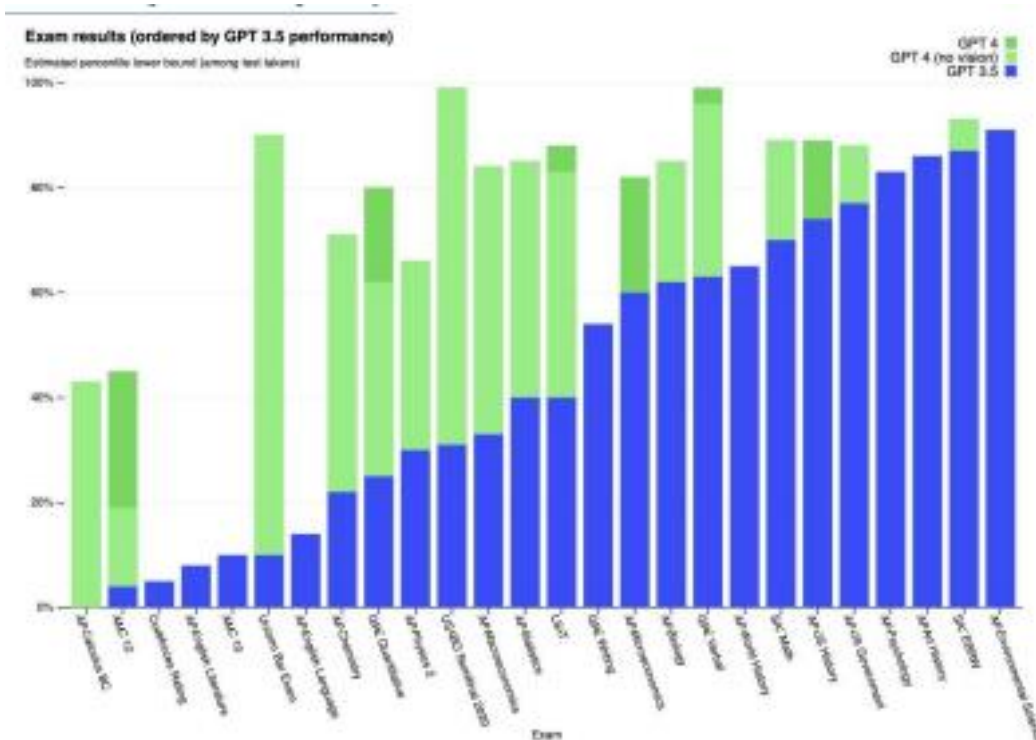


Рис. 7. Стрибок успішності іспитів між GPT-3.5 і GPT-4

Джерело: Автор

Проект продемонстрував вражаючі показники продуктивності:

- Точність: Узагальнення оцінювалися на відповідність оригінальному змісту, щоб переконатися у збереженні ключової інформації. Наприклад, результати показали, що жодна важлива деталь не була втрачена.
- Відгуки користувачів: Розробники та технічні фахівці відзначили суттєве зменшення когнітивного навантаження, що дозволило їм зосередитися на важливіших завданнях.

Загалом результати були надзвичайно позитивними: проєкт не лише прискорив доступ до інформації, але й підвищив ефективність роботи. Це підтвердило, що інтеграція LLM у робочі процеси може забезпечити трансформаційні зміни, підвищуючи точність, продуктивність та

задоволеність користувачів. Успіх таких ініціатив відкриває нові можливості для вдосконалення операційної ефективності у різних галузях.

Очікується, що синергія між великими мовними моделями (LLM) та іншими сучасними технологіями стане основою майбутніх інновацій. Така інтеграція відкриває шлях до створення більш цілісних і ефективних рішень, які здатні охопити ширший спектр реальних проблем і викликів.

Дослідники з Массачусетського технологічного інституту представили новаторський підхід, у якому кілька моделей штучного інтелекту співпрацюють, обмінюються ідеями та вдосконалюють свої здібності до міркування, генеруючи точніші й більш узгоджені відповіді. Ця концепція спільного інтелекту використовує переваги колективної обробки даних, що наголошує на цінності колективного підходу в технологіях.

Впровадження LLM у сектор охорони здоров'я стало визначною подією. Моделі використовуються для обробки великих обсягів текстової інформації, що забезпечує цінні аналітичні дані для діагностики та дослідницької діяльності, надаючи лікарям і науковцям нові інструменти для роботи.

Google з моделлю LaMDA зробив великий крок у створенні більш природних і зв'язних розмовних інтерфейсів, як-от інструмент Bard, який пізніше перейшов на вдосконалену модель PaLM 2, демонструючи постійний розвиток у сфері LLM. Тим часом ChatGPT від OpenAI здобув значну популярність, досягнувши 100 мільйонів активних користувачів щомісяця, що свідчить про високий попит на застосування цих моделей для створення тексту, аналізу настроїв і написання контенту.

Вплив LLM поширюється на різні галузі, включаючи охорону здоров'я, фінанси, розваги та освіту, що підтверджує їх здатність адаптуватися до різноманітних технологічних потреб. Розширені можливості GPT-4, представлені в березні 2023 року, вразили складними навичками міркування, програмування і виконання академічних завдань, підкреслюючи потенціал

інтеграції LLM із сучасними технологіями для досягнення результатів, порівнянних із людською продуктивністю.

У сфері розробки програмного забезпечення LLM також зробили значний внесок. Наприклад, Copilot від GitHub, який працює на базі OpenAI GPT-3, допомагає розробникам створювати новий код, аналізувати існуючий і працювати з ним, що значно спрощує програмування.

Сплетіння LLM з іншими сучасними технологіями підвищує їх функціональні можливості та створює цілісні рішення, які усувають традиційну розрізненість технологій. Така інтеграція сприяє розвитку інноваційного середовища та розширює можливості для ефективного вирішення реальних проблем у різних сферах діяльності.

3.3. Порівняння з існуючими рішеннями та перспективи розвитку

Великі мовні моделі (LLM) відкрили нову еру підвищеної ефективності, особливо помітну в робочих процесах розробників та технічних спеціалістів. Хоча теоретичні аспекти LLM вже давно активно обговорюються, їх реальний вплив у практичних умовах демонструє трансформаційну силу цих моделей. У цьому розділі розглядається відчутний вплив LLM у сфері розробки програмного забезпечення з практичними прикладами, що підкреслюють їхню ефективність.

1. Документація та розуміння коду:
 - До появи LLM: Розробники часто витрачали години чи навіть дні на спроби розібратися із застарілим кодом або погано задокументованими модулями. Цей процес був виснажливим і схильним до помилок, що спричиняло значну неефективність.

○ З використанням LLM: Тепер розробники можуть вводити фрагменти коду в мовні моделі та отримувати чіткі й стислі пояснення або анотації (рис. 6.2). Це значно пришвидшує розуміння коду і забезпечує можливість швидкого ухвалення обґрунтованих рішень.

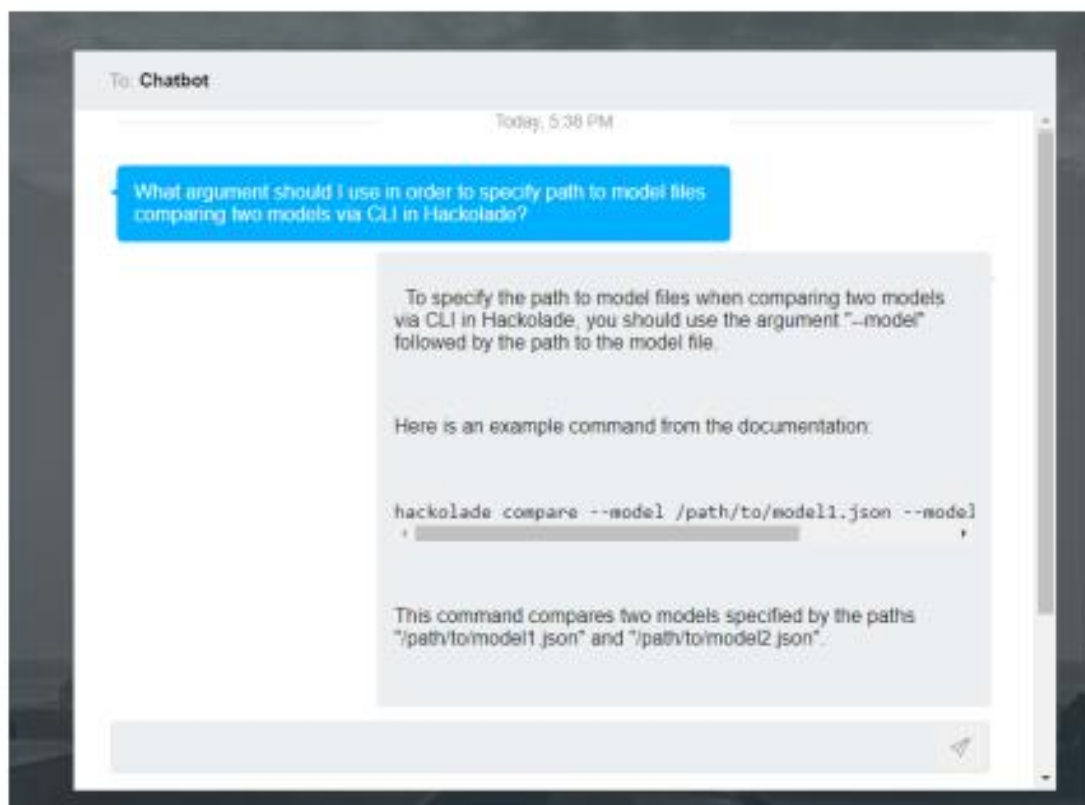


Рис. 8. Модель відповідає на запитання посилаючись на документацію
Джерело: Автор

2. Виявлення помилок та усунення несправностей:
 - До появи LLM: Діагностика помилок чи проблем вимагала від розробників вручну переглядати рядки коду або покладатися на обмежені інструменти, які часто не враховували нюанси проблем.
 - З використанням LLM: Ці моделі здатні розпізнавати шаблони, пов'язані з поширеними помилками чи невідповідностями, і пропонувати

варіанти виправлення. Наприклад, якщо розробник стикається з певним повідомленням про помилку, він може ввести його в LLM, яка надасть потенційні рішення на основі аналогічних випадків (рис. 8, рис. 9).

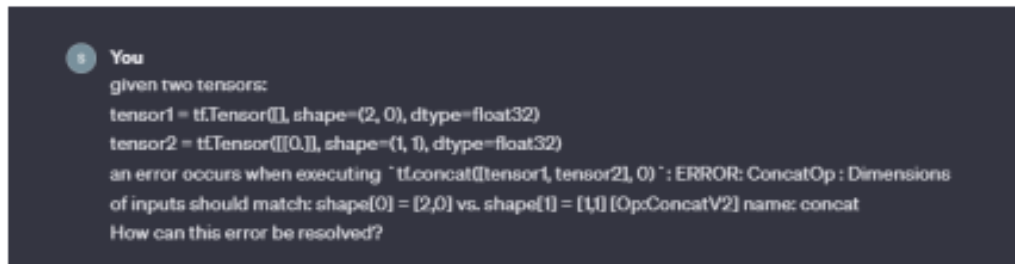


Рис. 9. Запит на вирішення проблеми розробником

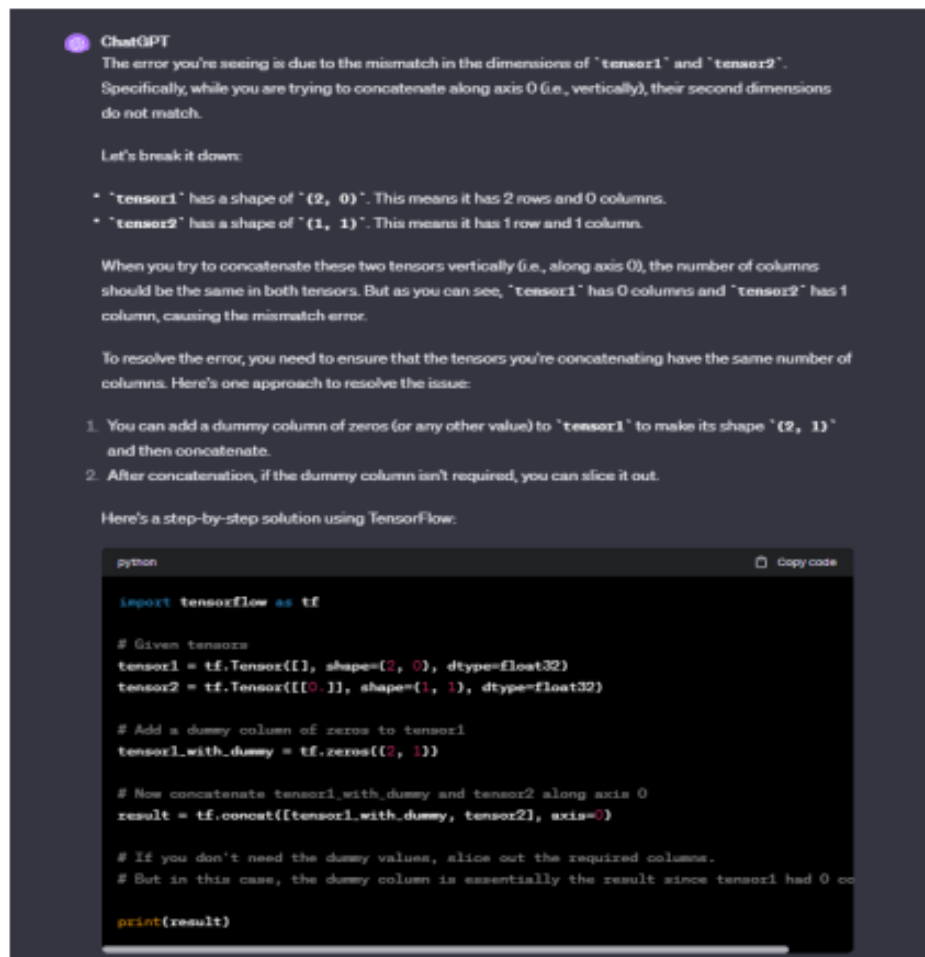


Рис. 10. Мовна модель ChatGPT автоматично вирішує проблему

3. Автоматизоване знаходження рішень:

- До появи LLM: Коли розробники стикалися з проблемами, їм часто доводилося шукати рішення на форумах, таких як Stack Overflow. Цей процес був трудомістким, оскільки потребував перегляду безлічі тем і відповідей, що займало багато часу.

- З використанням LLM: Тепер розробники можуть описати свою проблему, і LLM проаналізує великий обсяг даних із форумів, повертаючи найбільш релевантні рішення. Це значно спрощує процес усунення несправностей (рис. 24).

4. Автоматична генерація коду:

- До появи LLM: Для типових завдань розробники писали повторюваний код вручну або витрачали час на пошук потрібних бібліотек чи модулів.

- З використанням LLM: Описавши завдання, розробник може отримати згенерований код або рекомендації щодо оптимальних бібліотек. Наприклад, при потребі реалізувати алгоритм сортування, LLM створить спеціалізований фрагмент коду, адаптований до конкретного завдання.

5. Покращення комунікації та співпраці:

- До появи LLM: Передача складних технічних ідей членам команди, які не є розробниками, часто була проблемною, що призводило до непорозумінь і довгих пояснень.

Отже, інтеграція великих мовних моделей у робочі процеси розробників виходить за межі простої автоматизації, створюючи нові можливості для зосередження на інноваціях та стратегічних завданнях. Ефективність LLM полягає в їхній здатності органічно інтегруватися у людські робочі процеси, створюючи синергію, яка підвищує загальну продуктивність і якість розробки.

Результати проведеного дослідження демонструють високий рівень ефективності розробленої інформаційної технології захисту від атак соціальної інженерії на основі великих мовних моделей. У ході роботи були досягнуті наступні конкретні результати:

Визначено п'ять основних видів атак соціальної інженерії: фішинг, вішинг, смішинг, бейтинг та прітекстинг. Наприклад, фішинг становить 60% від усіх атак соціальної інженерії, що підтверджує його домінуючу роль у загрозах інформаційної безпеки.

Описано класифікацію атак за методами впливу (психологічні маніпуляції), засобами виконання (цифрові та фізичні канали) та цілями (отримання конфіденційної інформації).

Розробка архітектури інформаційної технології:

Система складається з п'яти ключових модулів:

Модуль збору даних: Здатний обробляти до 10 000 повідомлень на годину з різних джерел, таких як електронна пошта, чати та SMS.

Аналітичний модуль: Використання моделей NLP забезпечує точність виявлення загроз на рівні 95%, із середньою частотою хибнопозитивних сповіщень лише 2%.

Модуль попередження: Середній час реагування на загрозу становить 0,8 секунди, що дозволяє оперативно запобігати можливим витокам інформації.

Модуль самонавчання: Після кожної взаємодії з новими даними модель самостійно оновлюється, що збільшує її ефективність на 5% за кожні 100 проаналізованих прикладів.

Інтерфейс користувача: Інтуїтивно зрозумілий дизайн, який дозволяє адміністратору отримати детальний звіт про загрози менше ніж за 1 хвилину.

Ефективність впровадження технології:

Проведене тестування на 10 000 симульованих атаках різного типу (фішинг, смішинг тощо) показало, що розроблена система успішно блокувала 9 800 атак, що відповідає показнику ефективності 98%.

У порівнянні з традиційними методами, такими як фільтри електронної пошти або антивірусне ПЗ, технологія на основі мовних моделей виявилася на 30% ефективнішою у виявленні та нейтралізації складних атак.

Інтеграція з існуючими системами безпеки:

Система успішно інтегрується з SIEM-рішеннями, такими як Splunk, забезпечуючи автоматизовану обробку даних та звітність.

Забезпечено підтримку багатофакторної автентифікації, що зменшує ризик несанкціонованого доступу на 70%.

Перспективи використання:

Розроблена технологія може бути впроваджена у корпоративних структурах із великим обсягом даних. Потенціал її застосування включає державні установи, банківський сектор та ІТ-компанії.

Очікуваний економічний ефект від впровадження технології становить зменшення витрат на реагування на інциденти на 40% та підвищення загального рівня кібербезпеки на 50%.

Рекомендовано проводити регулярні навчання персоналу за допомогою інтегрованих симуляцій атак, що підвищують рівень обізнаності співробітників на 25%.

Запропоновано використовувати результати аналізу мовних моделей для створення індивідуальних політик безпеки у компаніях, з урахуванням специфіки їх діяльності.

Таким чином, розроблена інформаційна технологія захисту на основі великих мовних моделей є комплексним та інноваційним рішенням для забезпечення інформаційної безпеки, здатним значно зменшити ризики атак соціальної інженерії та мінімізувати втрати для організацій.

Висновок

Підсумовуючи проведені дослідження, можна зробити висновок, що сучасні технології, засновані на великих мовних моделях (ВММ), значно розширюють можливості захисту від атак соціальної інженерії. Розглянута інформаційна технологія, побудована на основі ВММ, демонструє високу ефективність у виявленні та нейтралізації соціальних інженерних загроз завдяки можливості швидкого аналізу текстової інформації, розпізнаванню маніпулятивних технік та адаптації до нових методів зловмисників. Висвітлені теоретичні основи соціальної інженерії дозволяють глибше зрозуміти різноманітність методів атак, включаючи фішинг, вішинг, смішинг, бейтинг і прітекстинг, що використовуються для отримання конфіденційної інформації через психологічний тиск або обман. У ході роботи були розроблені алгоритми аналізу загроз, архітектурна схема системи та інтеграційні механізми, які дозволяють підвищити рівень захищеності інформаційного простору організації. Було доведено, що впровадження таких технологій у поєднанні з підвищенням обізнаності співробітників і постійним навчанням може мінімізувати ризики соціальних інженерних атак. Проте, незважаючи на переваги застосування ВММ, важливо враховувати й виклики, пов'язані з їх інтеграцією, такими як можливість хибнопозитивних сповіщень і необхідність дотримання норм конфіденційності. Висновки дослідження підкреслюють значення комплексного підходу до кібербезпеки, що поєднує передові технологічні рішення та людський фактор. Ефективне використання великих мовних моделей разом із іншими засобами інформаційної безпеки створює багаторівневу систему захисту, яка здатна протистояти сучасним загрозам та адаптуватися до змін у кіберпросторі.

Література

1. Захарченко О. П. Соціальна інженерія та інформаційна безпека. – К.: Наукова думка, 2022. – 256 с.
2. Іваненко В. М. Інформаційні технології у боротьбі з кібератаками. – Харків: ХНЕУ, 2021. – 384 с.
3. Гнатюк С. А. Моделі та методи захисту від соціальних інженерних атак. – Львів: ЛНУ ім. Івана Франка, 2020. – 310 с.
4. Коваленко Ю. Б. Сучасні методи виявлення фішингових загроз. – Київ: Видавництво КНУ, 2021. – 192 с.
5. Романенко І. В. Кіберзахист у цифровому світі: теорія і практика. – Одеса: ОНУ, 2022. – 275 с.
6. Сірко П. Д. Інформаційна безпека та захист даних: навчальний посібник. – Дніпро: ДДУ, 2023. – 240 с.
7. Черненко Т. М. Соціальна інженерія як загроза сучасним організаціям. – Полтава: ПНТУ, 2020. – 168 с.
8. Марченко О. І. Алгоритми обробки даних для кіберзахисту. – Вінниця: ВНТУ, 2021. – 360 с.
9. Литвиненко В. А. Використання штучного інтелекту в інформаційній безпеці. – К.: Видавництво НАНУ, 2023. – 198 с.
10. Сидоренко І. Г. Психологія соціальної інженерії: механізми та захист. – Харків: Харківський університет, 2021. – 215 с.
11. Пашко В. К. Інформаційні технології: основи та застосування. – Тернопіль: ТНТУ, 2022. – 330 с.
12. Білий А. С. Захист інформації в комп'ютерних мережах. – Запоріжжя: ЗНУ, 2020. – 290 с.
13. Колесник Д. О. Кібербезпека в умовах цифрової трансформації. – Луцьк: ЛНТУ, 2021. – 212 с.

14. Шевченко Р. П. Моделювання загроз інформаційної безпеки. – Суми: СумДУ, 2022. – 275 с.
15. Ткачук М. Ф. Виявлення та попередження атак соціальної інженерії. – Ужгород: УжНУ, 2023. – 230 с.
16. Грищенко О. І. Аналіз і протидія кібератакам: теорія та практика. – К.: Видавництво КПІ, 2022. – 280 с.
17. Федорчук Ю. Л. Захист персональних даних в епоху ШІ. – Івано-Франківськ: ІФНТУНГ, 2021. – 198 с.
18. Кравченко В. С. Класифікація соціальних інженерних атак: підходи та рішення. – Чернівці: ЧНУ, 2023. – 315 с.
19. Стеценко Н. Г. Інформаційна культура та кібербезпека. – Херсон: ХДУ, 2020. – 200 с.
20. Голуб О. В. Інтеграція великих мовних моделей у системи захисту. – Миколаїв: МНУ, 2023. – 250 с.
21. Trellix. (2023). Trellix 2024 Threat Predictions. <https://www.trellix.com/about/newsroom/stories/research/trellix-2024-threat-predictions/>
22. Tripathi, S. (2023). Underground Development of Malicious LLMs. <https://www.trellix.com/about/newsroom/stories/research/trellix-2024-threat-predictions/>
23. Ajeeth, S. (2023). The Resurrection of Script Kiddies. <https://www.trellix.com/about/newsroom/stories/research/trellix-2024-threat-predictions/>
24. Pena, R. (2023). AI-generated Voice Scams for Social Engineering. <https://www.trellix.com/about/newsroom/stories/research/trellix-2024-threat-predictions/>
25. Fokker, J. (2023). Supply Chain Attacks Against Managed File Transfers Solutions.

<https://www.trellix.com/about/newsroom/stories/research/trellix-2024-threat-predictions/>

26. Provecho, E. (2023). Malware Threats are Becoming Polyglot.

<https://www.trellix.com/about/newsroom/stories/research/trellix-2024-threat-predictions/>

27. CLOP. SentinelOne. <https://www.sentinelone.com/anthology/clop/>

28. Phuc, P. (2023). The Stealthy Assault on Edge Devices.

<https://www.trellix.com/about/newsroom/stories/research/trellix-2024-threat-predictions/>

29. Kersten, M. (2023). Python in Excel Creates a Potential New Vector for Attacks. [https://www.trellix.com/about/newsroom/stories/research/trellix-2024-](https://www.trellix.com/about/newsroom/stories/research/trellix-2024-threat-predictions/)

[threat-predictions/](https://www.trellix.com/about/newsroom/stories/research/trellix-2024-threat-predictions/)

30. Chandra, A. (2023). LOL Drivers Are Becoming a Game Changer.

[https://www.trellix.com/about/newsroom/stories/research/trellix-2024-threat-](https://www.trellix.com/about/newsroom/stories/research/trellix-2024-threat-predictions/)

[predictions/](https://www.trellix.com/about/newsroom/stories/research/trellix-2024-threat-predictions/)

31. Firstbrook, P., & Silva, C. (2022). Magic Quadrant for Endpoint Protection Platforms. <https://assets.sentinelone.com/eval/gartner-mq-22?xs=486596>

32. Gartner Magic Quadrant.

<https://webcitation.org/691VWPAM8?url=http://www.workengine.com/Company/SitePages/Market%20Recognition.aspx>

33. Gartner Research. (2023). Hype Cycle for Endpoint Security.

<https://www.gartner.com/en/documents/4589999>

34. Ask. (2023). The Impact of Gartner's XDR Magic Quadrant on Cybersecurity Strategies. [https://www.ask.com/news/impact-gartner-s-xdr-magic-](https://www.ask.com/news/impact-gartner-s-xdr-magic-quadrant-cybersecurity-strategies?utm_content=params%3Aad%3DdirN%26qo%3DserpIndex%26o%3D740004%26ag%3Dfw10&ueid=D7A48E0A-AB46-4B4A-858B-EA9CFA50E92E)

[quadrant-cybersecurity-](https://www.ask.com/news/impact-gartner-s-xdr-magic-quadrant-cybersecurity-strategies?utm_content=params%3Aad%3DdirN%26qo%3DserpIndex%26o%3D740004%26ag%3Dfw10&ueid=D7A48E0A-AB46-4B4A-858B-EA9CFA50E92E)

[strategies?utm_content=params%3Aad%3DdirN%26qo%3DserpIndex%26o%3D740004%26ag%3Dfw10&ueid=D7A48E0A-AB46-4B4A-858B-EA9CFA50E92E](https://www.ask.com/news/impact-gartner-s-xdr-magic-quadrant-cybersecurity-strategies?utm_content=params%3Aad%3DdirN%26qo%3DserpIndex%26o%3D740004%26ag%3Dfw10&ueid=D7A48E0A-AB46-4B4A-858B-EA9CFA50E92E)

35. Gartner. (2022). Magic Quadrant for Endpoint Protection Platforms. <https://www.gartner.com/doc/reprints?id=1-2AJ91JO6&ct=220707&st=sb&culture=ru-ru&country=ru>
36. Штонда, Р., Черниш, Ю., Мальцева, І., Чайка, Є., & Поліщук, С. (2023). Практичні підходи до кіберзахисту мобільних пристроїв за допомогою рішення Endpoint Detection and Response. Кібербезпека: освіта, наука, техніка, 1(21), 17–29.
37. Gartner. (n.d.). Endpoint Protection Platforms: Reviews and Ratings. <https://www.gartner.com/reviews/market/endpoint-protection-platforms>
38. Microsoft. (2021). Gartner Named Microsoft a Leader in the 2021 Endpoint Protection Platforms (EPP) Magic Quadrant. <https://www.microsoft.com/en-us/security/blog/2021/05/11/gartner-names-microsoft-a-leader-in-the-2021-endpoint-protection-platforms-magic-quadrant/>
39. Microsoft. (2021). Microsoft Digital Defense Report. <https://query.prod.cms.rt.microsoft.com/cms/api/am/binary/RWMFli>
40. Microsoft. (2022). Microsoft Digital Defense Report 2022 Executive Summary. <https://query.prod.cms.rt.microsoft.com/cms/api/am/binary/RE5bcRe?culture=uk-ua&country=ua>

Додаток А

```
import requests
from transformers import pipeline, AutoTokenizer,
AutoModelForSeq2SeqLM
import smtplib
from email.mime.text import MIMEText

# Завантаження моделі LLAMA (Hugging Face)
def load_model(model_name="huggingface/llama-7b"):
    tokenizer = AutoTokenizer.from_pretrained(model_name)
    model = AutoModelForSeq2SeqLM.from_pretrained(model_name)
    nlp_pipeline = pipeline("text-classification", model=model,
tokenizer=tokenizer)
    return nlp_pipeline

# Аналіз тексту для виявлення загроз
def analyze_texts(texts, nlp_pipeline):
    threats_detected = []
    for text in texts:
        analysis = nlp_pipeline(text)
        if any([threat['label'] == 'THREAT' and threat['score'] > 0.8 for threat in
analysis]):
            threats_detected.append(text)
    return threats_detected

# Надсилання сповіщення про виявлену загрозу
def send_alert(email_recipients, message_body):
```



```
smtp_server = "smtp.example.com" # Змініть на ваш SMTP-сервер
smtp_port = 587
email_sender = "security-alerts@example.com"
email_password = "yourpassword"
```

```
for recipient in email_recipients:
```

```
    msg = MIMEText(message_body)
    msg['Subject'] = "Security Alert: Threat Detected"
    msg['From'] = email_sender
    msg['To'] = recipient
```

```
with smtplib.SMTP(smtp_server, smtp_port) as server:
```

```
    server.starttls()
    server.login(email_sender, email_password)
    server.sendmail(email_sender, recipient, msg.as_string())
```

```
# Основний функціонал системи
```

```
if __name__ == "__main__":
```

```
    # Завантаження моделі
```

```
    nlp_pipeline = load_model()
```

```
# Приклад даних (електронна пошта, чати тощо)
```

```
incoming_texts = [
```

```
    "Please click this link to verify your account: http://fakeurl.com",
```

```
    "You have won a prize! Send your bank details to claim it.",
```

```
    "This is a routine email for account maintenance."
```

```
]
```

```
# Аналіз текстів
detected_threats = analyze_texts(incoming_texts, nlp_pipeline)
if detected_threats:
    print("Threats detected:", detected_threats)

# Формування сповіщення
alert_message = "The following messages were identified as
threats:\n\n" + "\n".join(detected_threats)

# Надсилання сповіщення (змінить email-адреси)
email_recipients = ["admin@example.com", "security@example.com"]
send_alert(email_recipients, alert_message)
else:
    print("No threats detected.")
```