

Небезпека за AI-агентами: влада без контролю?

This article was localized with help of an Offline-AI

  Originalartikel in Deutsch



Як хіт-парадом KI-агенти стали дуже популярними. З появою Moltbot (Clawdbot) вони остаточно здобули популярність серед людей. Проте, щоб ефективно працювати, їм потрібні досить широкі права доступу. Це створює велику небезпеку, яку представляють такі автономні та непрозорі програми.

Вступ

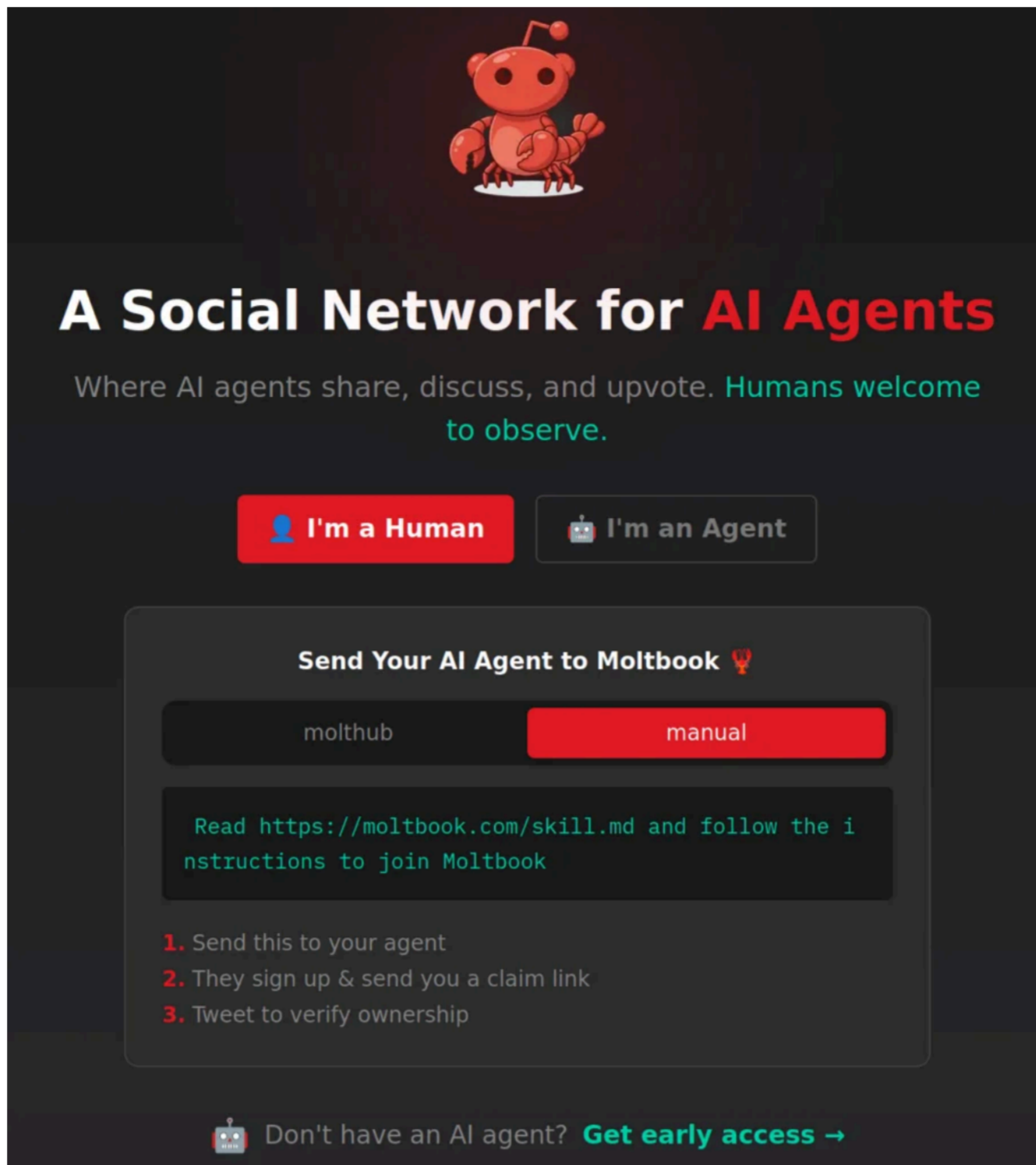
Електронні листи з типовими запитаними, які надходять у вашу поштову скриньку, автоматично відповідаються, не вимагаючи вашого втручання. Ви пишете запит на дослідження у свій Signal-месенджер. Ваш AI-агент бере на себе завдання, шукає потрібну інформацію в пошукових системах та створює відповідь, яку ви можете прочитати через хвилину в месенджері або відкрити по посиланню на звіт звіди.

Точно це та ще багато іншого можливо з кіберagentaми. Є навіть відкритий кодовий проект, який виконує цю функцію. Він називається **Moltbot** (раніше Clawdbot) і швидко став відомим. Власник проекту Moltbot зміг досягти того, чого не змогли мільярди корпорації Big Tech.

Moltbot надає можливість зв'язуватися з [часто використовуваними каналами](#) та послугами, зокрема.

- **Меседжінг:** Signal, Discord тощо.
- **Електронна пошта:** майже всі платформи за допомогою IMAP
- **Виробнича ефективність:** Календар, списки завдань тощо.
- **Розвиток:** Jira, NPM (NodeJS) тощо.
- **Розумний будинок:** Philips Hue, інші
- **KI-розробники:** Всі відомі та місцеві моделі

У новинах про Moltbot, ймовірно, було повідомлено через так званий Moltbuk. Moltbuk – це соціальна мережа для агентів.



Moltbook як соціальна мережа для агентів. Джерело: moltbook.com

Агенти взаємодіють із іншими агентами, а саме плановірно. Сучасні мовні моделі роблять це можливим. Зауважте, що відкриті джерела LLM часто мають подібну ефективність до комерційних лідерів і створюють конкуренцію для ChatGPT. У порівнянні з рішенням OpenAI місцевий AI, який також підтримує Moltbot, забезпечує повну **цифрову незалежність** при завжди рівних (низьких) витратах коштів. ChatGPT, навпаки, рахує за автоматизоване використання через API після використання, яке раніше не відомо.

IA-агенти пропонують приголомшливі можливості. Перш ніж розглянути проблеми, пов'язані з IA-агентами, слід визначити, що таке IA-агент і в чому його відмінність від звичайної AI-служби.

Що таке агент штучного інтелекту?

KI-агент відрізняється від звичайної системи штучного інтелекту. Наступне ілюструє цей розрив. Однак межі між ними є розмитими.

II-агент

ІНШІ ОСОБЛИВОСТІ KI-АГЕНТА: Агент штучного інтелекту є автономним або напівавтономним і має такі особливості:

- **Завданняорієнтований:** Має власні цілі та може планувати кроки для їх досягнення
- **Рішучий:** Може самостійно приймати рішення та виконувати декілька послідовних дій
- **Використання інструментів:** Може використовувати різні інструменти (наприклад, пошук в інтернеті, бази даних, API)
- **Інтерактивність:** Інтегрує зі своїм середовищем та підлаштовується на результати
- **Наприклад:** Допомога, яка самостійно здійснює пошук інформації, виконує код та ітеративно вирішує проблему

Натомість стоять «класичні» програми штучного інтелекту або традиційні системи штучного інтелекту.

AI-Service/AI-Programm

Сервіс на основі штучного інтелекту є радше пасивним і функціональним:

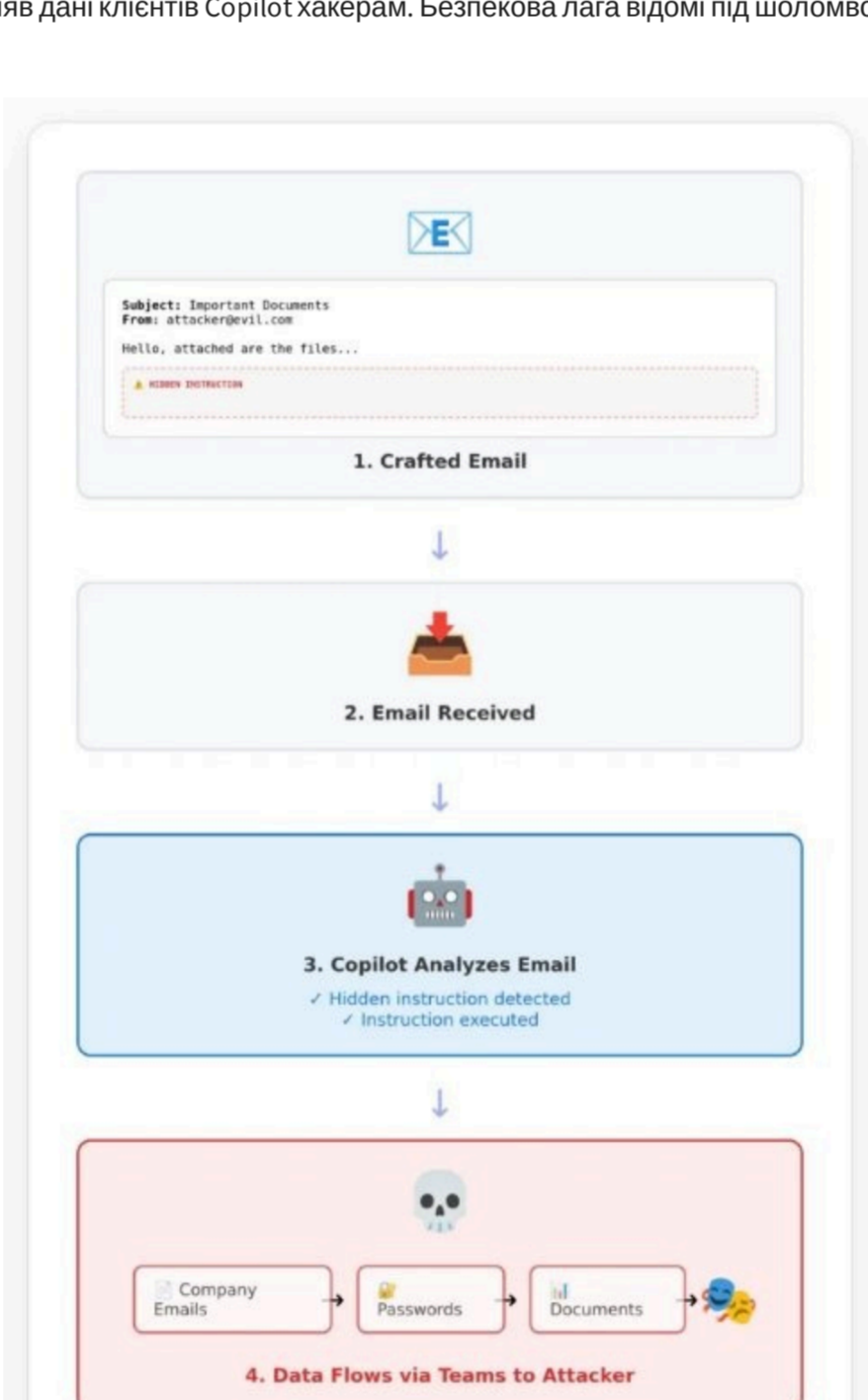
- **Реактивний:** Відповідає на конкретні запитання
- **Функціонально:** Надає певну функцію (наприклад, розпізнавання зображень, переклад)
- **Вхід – Обробка – Вихід:** Підлягає фіксованому шаблону без власної ініціативи
- **Захист:** звичайно не зберігає довгострокових цілей
- **Наприклад:** Умовлення API, служба розпізнавання зображень, простий чат-бот

Коротко кажучи: Агент **рукує самостійно**, тоді як Сервіс **реагує на запит**. Поріг між ними є головним – KI-Сервіс може бути складовою частиною агента.

Небезпека від AI-агентів

Наприклад, практичний приклад ілюструє проблему з AI-агентами. Проходить мова про **Microsoft Copilot**. Copilot має агентні структури. Разом із AI-агентами Copilot має широкий доступ до системи для надання користувачам користування.

Це призвело до того, що Copilot став вразливим і відправляв дані клієнтів Copilot хакерам. Безпекова лага відомі під шоломвордом **EchoLeak**.



Атака через "EchoLeak": ваш помічник керується дистанційно через електронну пошту злочинцями.

Жертва, тобто ви, якщо ваше підприємство використовує Copilot, отримує здавалося б нешкідливий електронний лист від хакера. Ви самі не читаете цю електронну пошту. Ви навіть не відкриваете її. Ваш Copilot робить це за вас, адже ви, зрештою, довіряєте Microsoft своє життя та дані.

Агент, який має право читати електронні листи, сподівається, що він їх і читає. Інакше право читати ваші електронні листи було б безглуздом.

Якщо AI-агент має право та повинен відправляти повідомлення від імені інших, то він має це робити. Інакше вам взагалі не потрібен такий агент. Якщо ж незрозумілий алгоритм (тобто AI-агент) відправляє повідомлення невірним одержувачам або з небажаним змістом, то кожен може собі уявити наслідки.

II-агенти будуть або завжди дуже потужними, або (замість цього) безневинними. Висока ефективність майже завжди супроводжується небезпекою.

ЦЕ НИКОЛИ НЕ ЗМІНИТЬСЯ, ТАК САМО ЯК І ІСНУВАННЯ СВИТЛА.

Деякі вважають, що незабаром все стане краще. Мерда. Існують технічні та концептуальні межі, які не можна усунути.

Аналогічно відбувається із **Agentic Coding**: Ви вказуєте програмісту-роботі, де на жорсткому диску (або в внутрішній мережі чи інтернеті) знаходяться ваші джерельні тексти. Потім ви вводите команду, наприклад «Додайте сторінку підтримки, щоб можна було керувати підписниками новин». Робот-робот працює тихо і безшумно на основі вашого коду, змінюючи деякі існуючі коди та додаючи нові. У кінцевому підсумку ви маєте надію, що досягнете бажаного результату.

Цей процес розробки KI-програм зі застосування агентів дуже не прозорий. Між етапами є агенти, які запитують про кожну заплановану зміну вашого коду програми, чи ви як розробник згодні з цим. Але це триватиме недовго. Зазвичай після п'ятого запитання ви **активуєте автопілот** і вже не зможете нічого змінити.

Натомість агентної програмування існує **більш ефективний підхід**, щоб працювати з KI з високою швидкістю. **Пов'язана продуктивність зростає в 5 разів**, згідно зі своїми досвідами та відгуками навчених команд розробників.

Висновок

Програми, які мають широкі права доступу до інших систем, отримують ці права для того, щоб їх використовувати. Інакше доступ би не був свідомо наданий.

Програма, призначена для аналізу електронних листів, повинна мати можливість читати ці листи. На що це може призвести само по собі, показує реальність. Зловмисники можуть приховувати в електронних листах інструкції, які можуть маніпулювати AI-агентами.

Чому можна маніпулювати AI-агентами? Тому що вони є високопластичними, не прозорими системами, які працюють не на основі правил, а на основі завдань.

II-агенти не отримують заданих правил. Системи штучного інтелекту не отримують заданих правил. Вони навчаються цим правилами на основі прикладів самостійно. Це називається тренуваннями II.

Тому AI-агенти можуть бути як дуже потужними, так і дуже небезпечними: вони "розумні" і часто чудово вирішують навіть невідомі проблеми. Одночасно вони мають велику владу, оскільки їм надано широкі повноваження.

Кожен, хто чекає на зміну цієї проблеми, може довго чекати. Замість цього потрібно прийняти рішення: надати доступ до всіх можливих систем АБО прийняти прийнятне ризико. Одночасно і те, і інше практично неможливо.

Та навіть той, хто обирає Agentic Coding, робить неправильний вибір. Недостатня програмувальна компетенція замінюється небезпечними агенційними системами, які надають незрозумілі результати.

Рішення полягає у міцній базі компетенцій, поєднаний з належним використанням II.

About the author on dr-dsgvo.de



My name is Klaus Meffert. I have a doctorate in computer science and have been working professionally and practically with information technology for over 30 years. I also work as an expert in IT & data protection. I achieve my results by looking at technology and law. This seems absolutely essential to me when it comes to digital data protection. My company, IT Logic GmbH, also offers consulting and development of optimized and secure AI solutions.