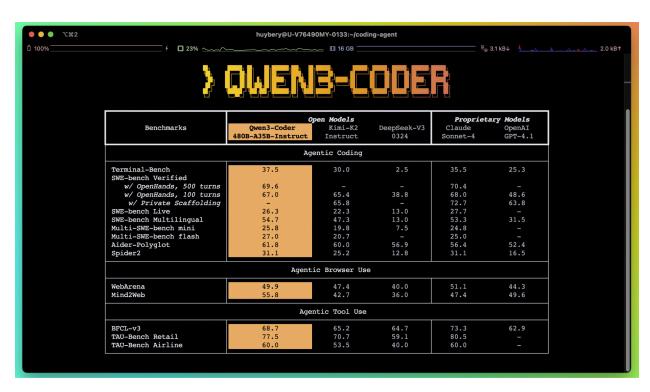
# Qwen3-Coder: Архітектурний та продуктивний аналіз передової агентної моделі коду





Розділ 1: Вступ та ключові висновки

1.1. Визначення нової парадигми: від автодоповнення коду до агентної інженерії

У ландшафті штучного інтелекту, що стрімко розвивається, поява Qwen3-Coder знаменує собою значний еволюційний стрибок. Ця модель, розроблена командою Qwen в Alibaba Cloud, виходить за межі традиційних асистентів кодування, які здебільшого зосереджені на генерації ізольованих фрагментів коду або функцій. Натомість Qwen3-Coder представляє нову категорію "агентних" моделей, спроектованих для автономного вирішення складних, багатоетапних завдань програмної інженерії. Цей звіт представляє глибокий аналіз архітектури, продуктивності та практичного застосування Qwen3-Coder, позиціонуючи її не просто як інструмент для написання коду, а як платформу для автоматизації комплексних робочих процесів розробки.

Фундаментальна відмінність полягає у переході від "генерації коду" до "автоматизації робочих процесів". Якщо попередні покоління моделей були ефективними у відповіді на конкретні запити, як-от "напиши функцію сортування", то Qwen3-Coder створена для виконання завдань на кшталт "проведи рефакторинг цього застарілого сервісу автентифікації, онови його залежності, перепиши тести та підготуй pull request". Ця здатність обробляти "невпорядковану, взаємопов'язану роботу, що визначає реальну програмну інженерію" <sup>4</sup>, є її ключовою ціннісною пропозицією. Така зміна фокусу вимагає не лише більшої обчислювальної потужності, а й принципово іншої архітектури та методології навчання, які дозволяють моделі розуміти контекст на рівні цілих репозиторіїв, планувати послідовність дій, використовувати зовнішні інструменти та виправляти власні помилки на основі зворотного зв'язку від середовища.

# 1.2. Ринкове позиціонування та конкурентне середовище

Qwen3-Coder виходить на ринок як потужний гравець у сегменті моделей з відкритим вихідним кодом, кидаючи виклик як пропрієтарним, так і іншим відкритим рішенням. Її продуктивність в агентних завданнях порівнюється з провідними закритими моделями, такими як Claude Sonnet від Anthropic. Водночас вона конкурує з найкращими відкритими моделями, призначеними для кодування, зокрема Kimi K2 та DeepSeek V3.

Однак ключовим стратегічним диференціатором Qwen3-Coder є її ліцензія. Модель випущена під дозвільною ліцензією Apache 2.0, що усуває значні бар'єри для комерційного використання та створення похідних продуктів. Че робить її особливо привабливою для підприємств, які прагнуть створювати власні, безпечні та кастомізовані інструменти для розробки на основі ШІ, не потрапляючи в залежність від постачальників пропрієтарних рішень. Поєднання передової продуктивності з відкритим доступом позиціонує Qwen3-Coder як потенційну фундаментальну технологію для наступного покоління корпоративних інструментів розробки, інтелектуальних агентів та автоматизованих конвеєрів розробки.

# 1.3. Ключові висновки звіту

Цей аналіз надає всебічну оцінку можливостей Qwen3-Coder, з якої можна зробити наступні ключові висновки:

- **Архітектурна перевага:** Поєднання архітектури Mixture-of-Experts (MoE) з величезним контекстним вікном створює потужну, але ефективну основу для розуміння коду в масштабах цілого репозиторію. МоЕ-дизайн дозволяє моделі мати величезну загальну кількість параметрів (480 мільярдів), активуючи лише невелику їх частину (35 мільярдів) для кожного токена, що забезпечує баланс між продуктивністю та обчислювальними витратами.<sup>2</sup>
- Профіль продуктивності: Модель демонструє найсучасніші результати на бенчмарках, що імітують реальні інженерні завдання, такі як SWE-Bench Verified, що підтверджує її основну силу в агентних сценаріях. Водночас її продуктивність на традиційних бенчмарках для генерації коду є висококонкурентною, але не завжди домінуючою, що вказує на спеціалізацію моделі. 4
- Практична корисність: Qwen3-Coder відмінно справляється зі складними, реальними інженерними робочими процесами, такими як модернізація застарілих систем, розробка функціоналу, що охоплює декілька сервісів, та налагодження розподілених систем. Це виходить далеко за рамки простої генерації окремих функцій.<sup>4</sup>
- **Екосистема з відкритим кодом:** Поєднання дозвільної ліцензії Apache 2.0 та потужних супутніх інструментів, зокрема інтерфейсу командного рядка Qwen Code, позиціонує Qwen3-Coder як фундаментальну модель для створення корпоративних, самостійно розміщуваних платформ розробки на основі ШІ.<sup>7</sup>

Основна ідея, що лежить в основі Qwen3-Coder, полягає у стратегічному зміщенні акценту з "генерації коду" на "автоматизацію робочих процесів". Філософія дизайну моделі віддає пріоритет здатності справлятися зі складною, взаємопов'язаною роботою, яка визначає сучасну програмну інженерію. Це стає можливим завдяки її архітектурі, зокрема величезному контекстному вікну, необхідному для аналізу цілих репозиторіїв <sup>4</sup>, та спеціалізованому навчанню на основі навчання з підкріпленням (RL) з довгим горизонтом, розробленому для багатоходових взаємодій. Показники продуктивності підтверджують цей фокус: модель досягає передових результатів на SWE-Bench, який симулює реальні проблеми в репозиторіях <sup>4</sup>, але демонструє більш змішані результати на простіших, однофункціональних бенчмарках. Таким чином, ціннісна пропозиція моделі полягає не просто в "написанні кращого коду", а в "автономному управлінні складними завданнями розробки" — фундаментально іншій та більш амбітній меті, що має значні наслідки для її оцінки та впровадження.

# Розділ 2: Архітектурний план Qwen3-Coder

#### 2.1. Перевага архітектури Mixture-of-Experts (MoE)

В основі виняткової продуктивності Qwen3-Coder лежить архітектура Mixture-of-Experts (MoE). Флагманська модель, Qwen3-Coder-480B-A35B-Instruct, має загальну кількість параметрів у 480 мільярдів, але під час обробки кожного токена активує лише 35 мільярдів параметрів. Цей підхід є ключовим для досягнення балансу між величезною ємністю моделі та ефективністю інференсу.

Система складається зі 160 окремих нейромережевих "експертів", кожен з яких потенційно спеціалізується на різних аспектах коду. Для обробки кожного вхідного токена механізм маршрутизації (gating network) динамічно обирає 8 найбільш релевантних експертів. Че дозволяє моделі використовувати знання значно більшої мережі, ніж це було б можливо з традиційною "щільною" архітектурою, при цьому зберігаючи обчислювальні витрати на керованому рівні. Така спеціалізація є надзвичайно важливою для домену кодування, оскільки різні експерти можуть навчатися розпізнавати та генерувати патерни для різних мов програмування (Python, Java, Rust), фреймворків (React, Django), парадигм (об'єктно-орієнтована, функціональна) або алгоритмічних підходів. Цей дизайн є вирішальним для досягнення високої продуктивності в реальних сценаріях розгортання.

# 2.2. Масштабування контексту: вікно в 1 мільйон токенів

Однією з найбільш визначних характеристик Qwen3-Coder є її здатність працювати з надзвичайно великими обсягами контексту. Модель нативно підтримує контекстне вікно довжиною 256,000 токенів, що вже є значним показником. Однак за допомогою методів екстраполяції, таких як YaRN (Yet another RoPE extensioN method), це вікно може бути розширене до 1 мільйона токенів.

Ця можливість підтримується передовими технічними рішеннями, інтегрованими в архітектуру. Технології, такі як Dual Chunk Attention (DCA) та MInference, дозволяють ефективно обробляти довгі послідовності. DCA розбиває довгий контекст на керовані частини, зберігаючи при цьому глобальну узгодженість, тоді як MInference є механізмом розрідженої уваги, що зменшує обчислювальне навантаження, фокусуючись лише на

критично важливих взаємодіях між токенами. 14

Стратегічне значення цієї характеристики неможливо переоцінити. Контекстне вікно такого розміру дозволяє моделі "утримувати в робочій пам'яті" цілі кодові бази, складні pull-реквести або великі обсяги документації. Че є необхідною передумовою для виконання справжніх агентних завдань на рівні репозиторію, таких як проведення наскрізного рефакторингу, аналіз залежностей між файлами або розуміння впливу змін в одному сервісі на інший. Без такої глибини контексту модель була б обмежена лише локальними, ізольованими змінами.

#### 2.3. Навчання в масштабі: дані, синтез та навчання з підкріпленням

Навчання Qwen3-Coder є багатоетапним та складним процесом, що відображає її спеціалізацію на агентних завданнях.

- Попереднє навчання (Pre-training): Основою моделі є попереднє навчання на масивному наборі даних об'ємом 7.5 трильйонів токенів. Критично важливим є те, що 70% цього набору складається з коду, що забезпечує глибоке розуміння синтаксису, семантики та патернів програмування. Важливою інновацією є використання попередньої моделі, Qwen2.5-Coder, для очищення, фільтрації та переписування "шумних" даних, що значно підвищує якість навчального корпусу. 2
- Пост-навчання (Code RL): Після попереднього навчання модель проходить етап доопрацювання за допомогою навчання з підкріпленням, керованого виконанням (execution-driven RL). Цей підхід базується на парадигмі "Важко вирішити, легко перевірити". Модель генерує код для вирішення завдань, а потім цей код виконується. Успішне виконання слугує позитивним сигналом для підкріплення, що безпосередньо підвищує відсоток успішного виконання коду, згенерованого моделлю.<sup>2</sup>
- Пост-навчання (Agent RL): Це ключовий етап, що наділяє модель агентними можливостями. Qwen3-Coder навчається за допомогою навчання з підкріпленням з довгим горизонтом (long-horizon RL). Для цього була створена масштабована система, здатна паралельно запускати 20,000 незалежних середовищ на інфраструктурі Alibaba Cloud. У цих середовищах модель вчиться виконувати багатоходові взаємодії, що включають планування, використання інструментів (наприклад, файлової системи або браузера), отримання зворотного зв'язку та прийняття рішень на його основі. Саме цей етап навчання дозволяє моделі виходити за рамки простої генерації коду та діяти як автономний агент.

Методологія навчання Qwen3-Coder демонструє замкнений цикл самовдосконалення. Замість пасивного навчання на статичному наборі даних, модель активно тренується в

динамічних середовищах, де вона вчиться на наслідках своїх дій. Використання попередньої моделі (Qwen2.5-Coder) для генерації синтетичних даних для її наступника (Qwen3-Coder) є першим кроком у цьому ітеративному циклі, що дозволяє створювати кращі навчальні матеріали. Наступним, більш динамічним етапом, є навчання з підкріпленням, де модель генерує дію (код або команду), середовище надає результат (успіх/невдача), і цей зворотний зв'язок використовується для оновлення політики моделі. Цей двосторонній підхід — покращення статичного набору даних, а потім динамічне тренування — є значно складнішим, ніж стандартне кероване доопрацювання. Це свідчить про те, що майбутнє навчання спеціалізованих моделей полягає у створенні таких масштабованих симуляційних середовищ, керованих зворотним зв'язком, що є значним інженерним викликом і конкурентною перевагою для команди Qwen.

# Розділ 3: Кількісний аналіз продуктивності

#### 3.1. Методологія бенчмаркінгу та застереження

Оцінка продуктивності Qwen3-Coder базується на наборі стандартизованих бенчмарків, кожен з яких вимірює різні аспекти можливостей моделі. Ключовими серед них є:

- **SWE-Bench:** Оцінює здатність моделі вирішувати реальні проблеми з репозиторіїв GitHub, що робить його надійним показником практичної інженерної компетентності.
- **LiveCodeBench:** Складається із завдань зі змагального програмування, що тестує алгоритмічні та оптимізаційні навички.
- HumanEval та MBPP (Mostly Basic Python Programming): Фундаментальні бенчмарки, що оцінюють функціональну коректність згенерованого коду на основі набору юніт-тестів.

Важливо зазначити, що при аналізі результатів виникають певні труднощі. У доступних джерелах спостерігаються розбіжності в показниках, а для деяких моделей відсутні дані по певних бенчмарках, що ускладнює пряме порівняння. Ця неоднозначність сама по собі є важливим висновком і буде відображена в аналізі.

# 3.2. Домінування в агентному кодуванні: SWE-Bench Verified

Найбільш переконливим доказом можливостей Qwen3-Coder є її продуктивність на бенчмарку SWE-Bench Verified. Модель демонструє найсучасніші результати серед моделей з відкритим вихідним кодом, причому ці результати були досягнуті без використання технік масштабування або аугментації під час тестування.

Значимість цього результату важко переоцінити. На відміну від синтетичних тестів, які часто фокусуються на ізольованих алгоритмічних задачах, SWE-Bench вимагає від моделі розуміння контексту існуючої кодової бази, виявлення причини проблеми та внесення точних змін для її вирішення. Успіх у цьому бенчмарку є прямим підтвердженням ефективності агентного підходу до навчання та архітектури з великим контекстним вікном, що дозволяє моделі ефективно працювати зі складними, реальними інженерними завданнями.

# 3.3. Продуктивність на традиційних та змагальних бенчмарках кодування

На інших, більш традиційних бенчмарках, картина є більш нюансованою.

- LiveCodeBench: На цьому бенчмарку, що складається із завдань зі змагального програмування, Qwen3-Coder показує себе як висококонкурентна модель. Однак вона не завжди займає перші місця, поступаючись деяким пропрієтарним моделям, таким як о4-mini від OpenAI. Показники для варіанту Qwen3-235В коливаються в різних джерелах: одне джерело вказує на результат 65.9% б, тоді як інше повідомляє про 70.7% для версії v5 бенчмарку 7, що підкреслює важливість врахування версії тесту.
- **HumanEval та MBPP:** Продуктивність на цих фундаментальних тестах є предметом певних розбіжностей. Деякі звіти стверджують, що Qwen3-Coder перевершує такі моделі, як GPT-4.1. Проте інші джерела описують її продуктивність як "сильну та конкурентну, але не завжди абсолютно найкращу", при цьому конкретні, перевірені числові показники часто відсутні в первинних джерелах. Ця відсутність чітких даних свідчить про те, що хоча модель є компетентною в генерації функціонально коректного коду, її основна сила та фокус розробки лежать в інших, більш складних агентних сценаріях.
- Якісні оцінки: Додаткові якісні тести показують, що Qwen3-Coder відмінно справляється зі стандартними завданнями середньої складності. Однак у деяких специфічних областях, таких як форматування складних візуалізацій в інтерфейсі користувача або вирішення завдань, що вимагають глибокого логічного міркування щодо нетипових патернів програмування, вона може поступатися таким моделям, як Кіті К2 або Claude 4.5

# 3.4. Порівняльний аналіз у таблицях

Для наочного представлення позиціонування Qwen3-Coder на ринку, нижче наведено дві зведені таблиці, що порівнюють її ключові характеристики та показники продуктивності з основними конкурентами.

#### Таблиця 1: Порівняльні специфікації моделей

Ця таблиця надає швидкий огляд ключових архітектурних та ліцензійних відмінностей, дозволяючи технічним керівникам оцінити масштаб, контекстні можливості та комерційну придатність кожної моделі.

| Модель               | Загальна<br>кількість<br>параметрів | Активні<br>параметри<br>(для МоЕ) | Макс.<br>довжина<br>контексту<br>(Нативна/Роз<br>ширена) | Ліцензія     |
|----------------------|-------------------------------------|-----------------------------------|----------------------------------------------------------|--------------|
| Qwen3-Coder<br>-480B | 480B                                | 35B                               | 256K / 1M                                                | Apache 2.0   |
| Qwen3-235B           | 235B                                | 22B                               | 256K / 1M                                                | Apache 2.0   |
| Claude 3<br>Opus     | Не<br>розголошуєть<br>ся            | Н/Д                               | 200K                                                     | Пропрієтарна |
| GPT-4o               | Не<br>розголошуєть<br>ся            | н/д                               | 128K                                                     | Пропрієтарна |
| DeepSeek-V3          | Не<br>розголошуєть<br>ся            | Не<br>розголошуєть<br>ся          | 128K                                                     | Пропрієтарна |
| Kimi K2<br>Instruct  | Не<br>розголошуєть<br>ся            | Н/Д                               | 200K                                                     | Пропрієтарна |

#### Таблиця 2: Зведена продуктивність на бенчмарках

Ця таблиця синтезує розрізнені дані про продуктивність в єдине представлення. Вона прямо відповідає на потребу в оцінці продуктивності, порівнюючи Qwen3-Coder з конкурентами на ключових бенчмарках. Важливо, що таблиця також вказує, де дані недоступні або суперечливі, надаючи чесну оцінку поточного стану знань.

| Модель               | SWE-Bench<br>Verified<br>(Pass@1) | LiveCodeBenc<br>h (Pass@1) | HumanEval<br>(Pass@1)              | MBPP<br>(Pass@1)                   |
|----------------------|-----------------------------------|----------------------------|------------------------------------|------------------------------------|
| Qwen3-Coder<br>-480B | ~69.6%                            | н/д                        | > GPT-4.1<br>(число не<br>вказано) | > GPT-4.1<br>(число не<br>вказано) |
| Qwen3-235B           | Н/Д                               | 65.9% - 70.7%              | Н/Д                                | н/Д                                |
| Claude<br>Sonnet 4   | ~65.2%                            | н/д                        | н/д                                | н/д                                |
| Claude 3<br>Opus     | н/д                               | 70.2%                      | н/д                                | 86.4%                              |
| GPT-4.1              | н/д                               | н/д                        | н/д                                | н/Д                                |
| DeepSeek-V3          | ~45-55%                           | Н/Д                        | н/д                                | 96.6%                              |
| Kimi K2<br>Instruct  | н/д                               | 70.4%                      | н/д                                | н/д                                |

Примітка: "Н/Д" (Немає даних) вказує на відсутність надійних числових показників у проаналізованих джерелах. Результати для HumanEval/MBPP для Qwen3-Coder описані якісно як такі, що перевершують GPT-4.1, але без конкретних чисел.

Джерела:  $^4$ 

# Розділ 4: Функціональні можливості в робочих

# процесах програмної інженерії

#### 4.1. Вирішення складних проблем та завдання системного рівня

Справжня цінність Qwen3-Coder розкривається не в ізольованих завданнях, а в її здатності вирішувати складні, системні проблеми, що є серцевиною сучасної розробки програмного забезпечення.

- Модернізація застарілих систем (Legacy System Modernization): Модель здатна виконувати комплексний аналіз застарілих кодових баз, виявляти вразливості безпеки, планувати та реалізовувати міграції на нові фреймворки або архітектурні патерни (наприклад, міграція з застарілих механізмів автентифікації на OAuth). Це включає здатність підтримувати зворотну сумісність під час внесення змін, що є критично важливим для мінімізації ризиків у виробничих системах.<sup>4</sup>
- Розробка крос-системного функціоналу (Cross-System Feature Development): Qwen3-Coder може керувати наскрізною реалізацією функціоналу, що охоплює декілька компонентів системи: від змін у бекенд-API та базах даних до відповідних оновлень у фронтенд-компонентах та конвеєрах розгортання. Вона може обробляти складну логіку, таку як впровадження обмеження частоти запитів (rate limiting), інтеграція платіжних систем або реалізація функціоналу для багатокористувацьких (multi-tenant) систем, що зачіпає кожен шар стеку.<sup>4</sup>
- Просунуте налагодження та аналіз першопричин (Advanced Debugging & Root Cause Analysis): Модель демонструє здатність розслідувати проблеми в розподілених системах. Вона може відстежувати поширення збоїв між мікросервісами, аналізувати взаємодії для виявлення вузьких місць у продуктивності та пропонувати системні виправлення, що усувають першопричину проблеми, а не лише її симптоми. 4

# 4.2. Автоматизація основних завдань розробки

Окрім системних завдань, Qwen3-Coder ефективно автоматизує повсякденні, але трудомісткі процеси розробки.

• Генерація та автодоповнення коду: Хоча в досліджених матеріалах бракує конкретних прикладів "автодоповнення" в стилі IDE, загальні можливості генерації

коду є надзвичайно потужними. Модель здатна не просто писати окремі функції, а генерувати цілі, готові до запуску програми. Наприклад, вона може за один запит створити повноцінну, грабельну версію гри "Тетріс", самостійно визначивши необхідні бібліотеки, структуруючи логіку гри та реалізувавши всі компоненти. <sup>20</sup> Це демонструє здатність до автономної розробки, а не простої трансляції запиту в код. <sup>9</sup>

- Виправлення помилок та рефакторинг: Можливості моделі в цій області добре ілюструє практичний приклад, де їй було доручено проаналізувати, виправити та покращити погано написаний Python-скрипт. Qwen3-Coder успішно ідентифікувала логічні помилки (наприклад, ризик ділення на нуль), потенційні помилки часу виконання (відсутність ключів у словнику) та погані практики (неефективні API-запити в циклі). Вона надала виправлену версію, а потім провела рефакторинг, додавши обробку помилок, покращивши імена змінних, впровадивши типізацію та розбивши код на логічні функції для кращої читабельності. Важливо відзначити і обмеження: модель виправила код, але не змінила фундаментально неефективний алгоритм (виконання окремих API-запитів замість пакетного), що вказує на її сильні сторони в покращенні якості коду, а не обов'язково в алгоритмічній оптимізації. 19
- Генерація Text-to-SQL: Хоча прямих прикладів для флагманської моделі Coder не знайдено, екосистема Qwen демонструє сильні можливості в цій сфері. Існують приклади успішного доопрацювання менших моделей Qwen для завдань Text-to-SQL <sup>23</sup>, а також використання старших моделей, як-от Qwen-Max, для створення чат-ботів, що взаємодіють з базами даних. <sup>25</sup> Крім того, Qwen3-Coder використовується в конвеєрах для генерації синтетичних наборів даних для навчання SQL-моделей. <sup>26</sup> Це дозволяє зробити висновок, що флагманська модель володіє потужними, хоча і неявно продемонстрованими, можливостями для перетворення запитів природною мовою на SQL.

# 4.3. Допоміжні можливості в розробці програмного забезпечення

Qwen3-Coder також автоматизує важливі, але часто другорядні завдання, що підвищує загальну продуктивність команди.

- **Автоматизована документація:** Модель відмінно справляється з генерацією вбудованих коментарів, файлів README, рядків документації (docstrings) та вичерпної документації для API. Ця функція є неоціненною для підтримки кодової бази в актуальному стані, спрощення процесу онбордингу нових членів команди та покращення загальної супроводжуваності проекту.<sup>7</sup>
- Обізнаність у питаннях безпеки: Хоча Qwen3-Coder не є формальним інструментом для статичного аналізу безпеки, вона може допомагати розробникам виявляти та запобігати поширеним вразливостям у коді. Її можна використовувати для перегляду ризикованих патернів, пропозицій щодо оновлення небезпечних

залежностей та впровадження найкращих практик безпеки в усьому стеку технологій, що робить її цінним асистентом під час код-рев'ю.<sup>7</sup>

Функціональні можливості Qwen3-Coder чітко вказують на те, що її цінність полягає не в заміні розробника для виконання одного завдання (наприклад, написання функції), а в розширенні його можливостей протягом усього робочого циклу. Її здібності безпосередньо спрямовані на найбільш трудомісткі та когнітивно навантажені аспекти сучасної програмної інженерії: роботу зі складністю, застарілим кодом та розподіленими системами. Набір функцій моделі націлений на ці "зони високого тертя" в розробці. Отже, її вплив слід вимірювати не кількістю згенерованих рядків коду за хвилину, а здатністю зменшити когнітивне навантаження на розробників, прискорити адаптацію до нових кодових баз та скоротити час, що витрачається на налагодження та супровід. Вона функціонує як "системний партнер для мислення", а не як простий "друкар коду".

# Розділ 5: Екосистема розробника та посібник з впровадження

#### 5.1. Варіанти доступу та розгортання

Для розробників існує кілька способів отримати доступ до Qwen3-Coder, що забезпечує гнучкість залежно від потреб у продуктивності, приватності та вартості.

- Доступ через API: Найпростіший спосіб почати роботу— це використання керованих кінцевих точок, що надаються такими платформами, як Alibaba Cloud Model Studio, OpenRouter та Together AI. Це усуває необхідність керувати інфраструктурою та дозволяє швидко інтегрувати модель у існуючі додатки.
- Локальне розгортання: Для сценаріїв, що вимагають максимальної приватності, роботи в офлайн-режимі або глибокої кастомізації, Qwen3-Coder можна розгортати локально. Інструменти, такі як Ollama, LMStudio та llama.cpp, підтримують моделі Qwen, дозволяючи запускати їх на власному обладнанні.<sup>8</sup>
- **Квантовані моделі:** Щоб зробити велику модель доступною для обладнання споживчого класу, спільнота та розробники Qwen надають квантовані версії (наприклад, FP8, GGUF). Ці моделі мають значно менший обсяг пам'яті, що дозволяє запускати їх на системах з обмеженими ресурсами, хоча і з певним компромісом у точності. <sup>13</sup>

# 5.2. Інтерфейс командного рядка (CLI) Qwen Code

Центральним елементом екосистеми Qwen3-Coder  $\varepsilon$  Qwen Code — інструмент командного рядка з відкритим вихідним кодом. Він  $\varepsilon$  форком Gemini CLI від Google, спеціально адаптованим для розкриття повного потенціалу моделей Qwen.<sup>2</sup>

Qwen Code оснащений покращеними парсерами та кастомізованими протоколами виклику функцій, які є необхідними для реалізації агентних робочих процесів безпосередньо з терміналу. Це дозволяє моделі взаємодіяти з локальною файловою системою, виконувати команди та автоматизувати складні багатоетапні завдання. На практиці розробники можуть використовувати Qwen Code для таких завдань, як отримання короткого огляду архітектури проекту, автоматична генерація повноцінних веб-додатків або виконання складних операцій з репозиторієм, таких як rebase або обробка pull-реквестів.<sup>2</sup>

# 5.3. Інтеграція з бібліотеками та IDE

Qwen3-Coder легко інтегрується в існуючі інструменти та робочі процеси розробників.

• Hugging Face Transformers: Модель повністю підтримується популярною бібліотекою transformers. Для її використання необхідно завантажити модель та токенізатор з Hugging Face Hub. Важливою технічною деталлю є вимога використовувати останню версію бібліотеки (transformers>=4.51.0), оскільки старіші версії можуть викликати помилку KeyError: 'qwen3\_moe' через відсутність підтримки архітектури МоЕ, що використовується в Qwen3.9 Приклад коду для генерації тексту:

Python

from transformers import AutoModelForCausalLM, AutoTokenizer

```
tokenizer = AutoTokenizer.from_pretrained(model_name)
model = AutoModelForCausalLM.from_pretrained(
    model_name,
    torch_dtype="auto",
    device_map="auto"
```

model name = "Qwen/Qwen3-Coder-480B-A35B-Instruct"

• **Інтеграція з IDE:** Qwen3-Coder може бути інтегрована в популярні середовища розробки, такі як Visual Studio Code та IDE від JetBrains, через розширення та плагіни. Крім того, інструменти, як-от Aider, дозволяють використовувати модель для парного програмування зі штучним інтелектом безпосередньо в терміналі, що забезпечує тісну інтеграцію з робочим процесом розробника.<sup>7</sup>

# 5.4. Найкращі практики для інференсу та промптингу

Для досягнення оптимальних результатів при роботі з Qwen3-Coder рекомендується дотримуватися наступних практик:

- Параметри семплювання: Рекомендовані налаштування для генерації включають temperature=0.7, top\_p=0.8, top\_k=20 та repetition\_penalty=1.05. Ці параметри допомагають збалансувати креативність та узгодженість відповідей. 9
- **Керування ресурсами:** При роботі з великими контекстами важливо встановити достатню максимальну довжину виводу (наприклад, 65,536 токенів), щоб модель не обрізала відповідь. На обладнанні з обмеженими ресурсами, у разі виникнення помилок браку пам'яті (Out-of-Memory, OOM), рекомендується зменшити довжину контексту до меншого значення, наприклад, 32,768 токенів.
- **Режим роботи:** Важливо знати, що інструкційні моделі Qwen3-Coder (-Instruct) працюють виключно в "не-мислячому" режимі. Це означає, що вони не генерують

спеціальні блоки <think></think> у своїх відповідях, тому вказувати параметр enable thinking=False більше не потрібно. $^1$ 

# Розділ 6: Стратегічні наслідки та майбутні перспективи

## 6.1. Вплив ліцензії Apache 2.0

Вибір ліцензії Apache 2.0 для Qwen3-Coder є стратегічним кроком з далекосяжними наслідками. Чля дозвільна ліцензія надає користувачам безстрокові, всесвітні, неексклюзивні та безоплатні права на використання, відтворення, розповсюдження та створення похідних робіт. Че робить модель безпечною для комерційного використання, що є ключовою перевагою порівняно з моделями з більш обмежувальними ліцензіями.

Для бізнесу це означає можливість створювати власні, внутрішні інструменти розробки на основі передової моделі ШІ, зберігаючи при цьому повний контроль над своїми даними та інфраструктурою. Це усуває ризики, пов'язані з залежністю від пропрієтарних постачальників, конфіденційністю коду та непередбачуваними змінами в ціноутворенні. Таким чином, Apache 2.0 позиціонує Qwen3-Coder не просто як потужний інструмент, а як фундаментальну, дружню до підприємств платформу, що сприяє інноваціям та демократизує доступ до передових агентних технологій.<sup>8</sup>

# 6.2. Поточні обмеження та відгуки спільноти

Незважаючи на вражаючі можливості, Qwen3-Coder має певні обмеження. Об'єктивний аналіз показує, що модель іноді стикається з труднощами при вирішенні завдань, що вимагають глибокого логічного міркування щодо нетипових або езотеричних патернів програмування, як-от просунуті техніки звуження типів у TypeScript. Також відзначаються проблеми з форматуванням складних елементів інтерфейсу користувача.<sup>5</sup>

Відгуки від спільноти розробників загалом позитивні, але також вказують на нюанси.

Деякі користувачі зазначають, що для певних відкритих та творчих завдань рівень міркувань моделі можна порівняти з рівнем "молодшого студента", тоді як провідні пропрієтарні моделі демонструють більш зрілий підхід. За Це свідчить про те, що хоча Qwen3-Coder є надзвичайно потужною у своїй спеціалізованій ніші агентного кодування, вона може не бути універсально найкращим вибором для абсолютно всіх типів завдань програмування.

#### 6.3. Майбутня дорожня карта та заключні зауваження

Команда Qwen продовжує активно працювати над вдосконаленням своїх моделей. Згідно з їхніми заявами, майбутні плани включають випуск моделей інших розмірів для зниження вартості розгортання, а також дослідження захоплюючого напрямку самовдосконалюваних агентів.<sup>2</sup>

На завершення, Qwen3-Coder є знаковою віхою в розвитку штучного інтелекту з відкритим вихідним кодом. Вона переконливо демонструє, що агентні можливості, які раніше були прерогативою провідних пропрієтарних моделей, можуть бути досягнуті та демократизовані. Модель встановлює новий стандарт для того, чого можна очікувати від відкритих рішень, зміщуючи фокус з простої генерації коду на автоматизацію складних інженерних робочих процесів. Її справжній вплив буде вимірюватися не лише показниками на бенчмарках, а й багатством екосистеми інструментів, додатків та інновацій, які будуть побудовані на її потужній та відкритій основі.

#### Джерела:

- 1. QwenLM/Qwen3-Coder: Qwen3-Coder is the code version ... GitHub, accessed August 23, 2025, <a href="https://github.com/QwenLM/Qwen3-Coder">https://github.com/QwenLM/Qwen3-Coder</a>
- 2. Qwen3-Coder: Agentic Coding in the World | Qwen, accessed August 23, 2025, <a href="https://qwenlm.github.io/blog/qwen3-coder/">https://qwenlm.github.io/blog/qwen3-coder/</a>
- 3. What You Need to Know About Qwen 3 Coder Today Vision Computer Solutions, accessed August 23, 2025, <a href="https://www.vcsolutions.com/blog/discover-qwen-3-coder-what-you-need-to-k-now-today/">https://www.vcsolutions.com/blog/discover-qwen-3-coder-what-you-need-to-k-now-today/</a>
- 4. Qwen3-Coder: The Most Capable Agentic Coding Model Now ..., accessed August 23, 2025, <a href="https://www.together.ai/blog/qwen-3-coder">https://www.together.ai/blog/qwen-3-coder</a>
- 5. Qwen3 Coder Performance Evaluation: A Comparative Analysis Against Leading Models, accessed August 23, 2025, https://eval.16x.engineer/blog/qwen3-coder-evaluation-results
- 6. Qwen3-Coder: The best Agentic Code AI, beats Kimi-K2 | by Mehul ..., accessed August 23, 2025, <a href="https://medium.com/data-science-in-your-pocket/qwen3-coder-the-best-agenti">https://medium.com/data-science-in-your-pocket/qwen3-coder-the-best-agenti</a>

- c-code-ai-beats-kimi-k2-1f8e6472c42b
- Qwen3 Coder: The Open-Source Al Coding Model Redefining Code Generation |
   Data Science Dojo, accessed August 23, 2025,
   <a href="https://datasciencedojo.com/blog/qwen3-coder/">https://datasciencedojo.com/blog/qwen3-coder/</a>
- 8. Introducing **Qwen-Code**: Alibaba's Open-Source CLI for Agentic Coding with Qwen3-Coder NYU Shanghai RITS, accessed August 23, 2025, <a href="https://rits.shanghai.nyu.edu/ai/introducing-qwen-code-alibabas-open%E2%80%91source-cli-for-agentic-coding-with-gwen3%E2%80%91coder/">https://rits.shanghai.nyu.edu/ai/introducing-qwen-code-alibabas-open%E2%80%91source-cli-for-agentic-coding-with-gwen3%E2%80%91coder/</a>
- 9. Qwen/Qwen3-Coder-480B-A35B-Instruct Hugging Face, accessed August 23, 2025, <a href="https://huggingface.co/Qwen/Qwen3-Coder-480B-A35B-Instruct">https://huggingface.co/Qwen/Qwen3-Coder-480B-A35B-Instruct</a>
- 10. Qwen3-Coder: Performance, Architecture & Access Zenn, accessed August 23, 2025, https://zenn.dev/saan/articles/68f1076b82d41e
- 11. qwen3-coder Ollama, accessed August 23, 2025, https://ollama.com/library/gwen3-coder
- 12. Qwen3 Coder API, Providers, Stats OpenRouter, accessed August 23, 2025, <a href="https://openrouter.ai/qwen/qwen3-coder">https://openrouter.ai/qwen/qwen3-coder</a>
- 13. Qwen3 Coder API Service Vertex AI Google Cloud console, accessed August 23, 2025, <a href="https://console.cloud.google.com/vertex-ai/publishers/qwen/model-garden/qwen3-coder-480b-a35b-instruct-maas?hl=ja">https://console.cloud.google.com/vertex-ai/publishers/qwen/model-garden/qwen3-coder-480b-a35b-instruct-maas?hl=ja</a>
- 14. Qwen/Qwen3-235B-A22B-Instruct-2507 Hugging Face, accessed August 23, 2025, https://huggingface.co/Qwen/Qwen3-235B-A22B-Instruct-2507
- 15. LiveCodeBench Benchmark Vals Al, accessed August 23, 2025, https://www.vals.ai/benchmarks/lcb-07-22-2025
- 16. LiveCodeBench Leaderboard Holistic and Contamination Free Evaluation, accessed August 23, 2025, <a href="https://livecodebench.github.io/leaderboard.html">https://livecodebench.github.io/leaderboard.html</a>
- 17. Qwen3 Technical Report, accessed August 23, 2025, https://arxiv.org/pdf/2505.09388
- 18. MBPP Benchmark (Code Generation) Papers With Code, accessed August 23, 2025, <a href="https://paperswithcode.com/sota/code-generation-on-mbpp">https://paperswithcode.com/sota/code-generation-on-mbpp</a>
- 19. Getting Started with Qwen3-Coder Analytics Vidhya, accessed August 23, 2025, <a href="https://www.analyticsvidhya.com/blog/2025/07/getting-started-with-qwen3-coder/">https://www.analyticsvidhya.com/blog/2025/07/getting-started-with-qwen3-coder/</a>
- 20. Hands-on Tutorial: Build Your Own Coding Copilot with Qwen3-Coder, Qwen Code, and Code Context Milvus, accessed August 23, 2025, <a href="https://milvus.io/blog/hands-on-tutorial-build-your-own-coding-copilot-with-qwen3-coder-qwen-code-and-code-context.md">https://milvus.io/blog/hands-on-tutorial-build-your-own-coding-copilot-with-qwen3-coder-qwen-code-and-code-context.md</a>
- 21. Qwen3 Coder COMPLETE Agentic Coding Test (Full Production App Build) YouTube, accessed August 23, 2025, https://www.youtube.com/watch?v=QPl8li1p-zY
- 22. Vibe Coded with Qwen 3 Coder in <1 hour : r/LocalLLaMA Reddit, accessed August 23, 2025, <a href="https://www.reddit.com/r/LocalLLaMA/comments/1m7u02i/vibe\_coded\_with\_qwen3 coder\_in1 hour/">https://www.reddit.com/r/LocalLLaMA/comments/1m7u02i/vibe\_coded\_with\_qwen3 coder\_in1 hour/</a>
- 23. fahmiaziz/qwen3-1.7B-text2sql Hugging Face, accessed August 23, 2025,

- https://huggingface.co/fahmiaziz/gwen3-1.7B-text2sgl
- 24. Ellbendls/Qwen-2.5-3b-Text\_to\_SQL Hugging Face, accessed August 23, 2025, <a href="https://huggingface.co/Ellbendls/Qwen-2.5-3b-Text\_to\_SQL">https://huggingface.co/Ellbendls/Qwen-2.5-3b-Text\_to\_SQL</a>
- 25. Generating SQL Queries with Alibaba Cloud's Qwen, accessed August 23, 2025, <a href="https://www.alibabacloud.com/blog/generating-sql-queries-with-alibaba-clouds-qwen">https://www.alibabacloud.com/blog/generating-sql-queries-with-alibaba-clouds-qwen</a> 602329
- 26. Text to sql Fireworks Al Docs, accessed August 23, 2025, <a href="https://docs.fireworks.ai/examples/text-to-sql">https://docs.fireworks.ai/examples/text-to-sql</a>
- 27. Qwen3 Coder (free) API, Providers, Stats OpenRouter, accessed August 23, 2025, https://openrouter.ai/gwen/gwen3-coder:free
- 28. New code benchmark puts Qwen 3 Coder at the top of the open models Reddit, accessed August 23, 2025,
  <a href="https://www.reddit.com/r/LocalLLaMA/comments/1mto8fa/new\_code\_benchmark">https://www.reddit.com/r/LocalLLaMA/comments/1mto8fa/new\_code\_benchmark</a>
  puts gwen 3 coder at the top/
- 29. Help: Qwen3-Coder + LM Studio + Continue.dev (VSCode) + Mac 64GB M3 Max 500 Internal Server Error, Even After Unsloth Fix Reddit, accessed August 23, 2025, <a href="https://www.reddit.com/r/LocalLLaMA/comments/1mf0fgj/help\_qwen3coder\_lm\_studio\_continuedev\_vscode\_mac/">https://www.reddit.com/r/LocalLLaMA/comments/1mf0fgj/help\_qwen3coder\_lm\_studio\_continuedev\_vscode\_mac/</a>
- 30. LICENSE · Qwen/Qwen3-Coder-480B-A35B-Instruct-FP8 at main Hugging Face, accessed August 23, 2025, <a href="https://huggingface.co/Qwen/Qwen3-Coder-480B-A35B-Instruct-FP8/blob/main/LICENSE">https://huggingface.co/Qwen/Qwen3-Coder-480B-A35B-Instruct-FP8/blob/main/LICENSE</a>
- 31. qwen3-coder:480b/license Ollama, accessed August 23, 2025, <a href="https://ollama.com/library/qwen3-coder:480b/blobs/d18a5cc71b84">https://ollama.com/library/qwen3-coder:480b/blobs/d18a5cc71b84</a>
- 32. Qwen 3 Benchmarks, Comparisons, Model Specifications, and More DEV Community, accessed August 23, 2025, <a href="https://dev.to/best\_codes/qwen-3-benchmarks-comparisons-model-specifications-and-more-4hoa">https://dev.to/best\_codes/qwen-3-benchmarks-comparisons-model-specifications-and-more-4hoa</a>