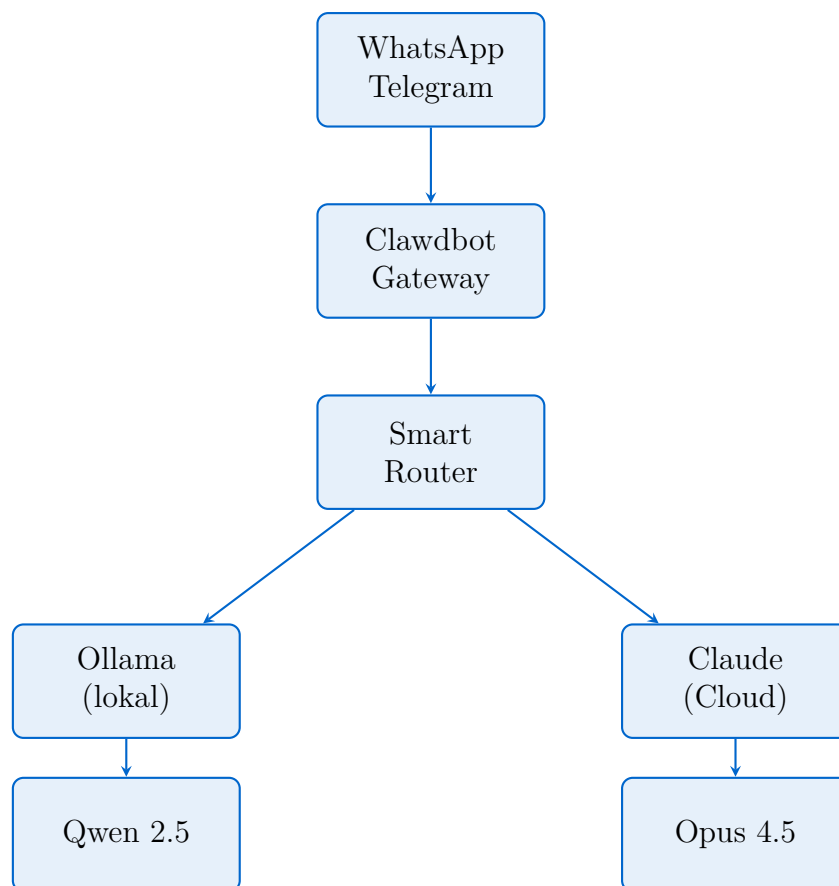


Ollama + Clawdbot

Hybrides LLM-System

Installation, Konfiguration & Best Practices



Version 1.0.0

27. Januar 2026

Dr. Jo

Inhaltsverzeichnis

1	Übersicht	5
1.1	Was ist Ollama?	5
1.2	Warum ein Hybrid-System?	5
1.2.1	Anwendungsfälle	5
1.3	Architektur	6
2	Voraussetzungen	7
2.1	Hardware	7
2.2	Software	7
2.2.1	Voraussetzungen prüfen	7
3	Installation	9
3.1	Ollama installieren	9
3.1.1	Via Homebrew (empfohlen)	9
3.1.2	Alternativer Download	9
3.2	Ollama Service starten	9
3.3	Qwen Modell herunterladen	10
3.3.1	Modell-Auswahl	10
3.3.2	Download-Befehle	10
3.4	Installation verifizieren	10
4	Konfiguration	11
4.1	Clawdbot konfigurieren	11
4.2	Hybride LLM-Konfiguration	11
4.3	Gateway neu starten	12
5	Qwen Plugin	13
5.1	Plugin erstellen	13
5.2	SKILL.md erstellen	13
5.3	package.json erstellen	13
5.4	Skill aktivieren	14
6	Testen	15
6.1	Ollama direkt testen	15
6.2	Interaktiver Test	15
6.3	Clawdbot-Integration testen	15
6.3.1	Via WhatsApp	15
7	Troubleshooting	17
7.1	Connection refused	17

7.2	Model not found	17
7.3	Ollama ist langsam	18
7.4	Routing funktioniert nicht	18
8	Erweiterte Nutzung	19
8.1	Weitere Modelle	19
8.1.1	Llama 3.1 (Meta)	19
8.1.2	Mistral (MistralAI)	19
8.1.3	CodeLlama (Code-spezialisiert)	19
8.2	Multi-Modell-Setup	19
8.3	Performance-Optimierung	20
8.3.1	Ollama-Konfiguration	20
8.3.2	Parameter-Erklärung	20
9	Sicherheit & Privacy	21
9.1	Warum Ollama sicherer ist	21
9.2	Best Practices	21
9.3	Logs löschen	21
A	Befehls-Referenz	23
A.1	Ollama-Befehle	23
A.2	Clawdbot-Befehle	23
B	Installations-Checklist	25
B.1	Vor der Installation	25
B.2	Installation	25
B.3	Tests	25

Kapitel 1

Übersicht

1.1 Was ist Ollama?

Ollama ist eine lokale LLM-Runtime, die große Sprachmodelle wie Qwen, Llama und Mistral direkt auf deinem Mac ausführt – komplett **offline** und **kostenlos**.

Hauptvorteile

- **Offline-Betrieb:** Funktioniert ohne Internetverbindung
- **Kostenlos:** Keine API-Gebühren
- **Privacy:** Daten verlassen niemals deinen Mac
- **Schnell:** Keine Netzwerk-Latenz

1.2 Warum ein Hybrid-System?

Die Kombination von **Ollama** (lokal) und **Claude API** (Cloud) bietet das Beste aus beiden Welten:

Aspekt	Ollama (lokal)	Claude (Cloud)
Kosten	Kostenlos	Pay-per-Use
Geschwindigkeit	Sehr schnell	Latenz durch Netzwerk
Privacy	Maximale Privacy	Daten in der Cloud
Qualität	Gut (14B Modell)	Exzellent (Opus)
Offline	Ja	Nein
Komplexität	Einfache Aufgaben	Komplexe Analysen

Tabelle 1.1: Vergleich Ollama vs. Claude

1.2.1 Anwendungsfälle

Ollama nutzen für:

- WhatsApp-Nachrichten beantworten

- Einfache Übersetzungen
- Textzusammenfassungen
- Schnelle Berechnungen
- Sensible Daten (Privacy)

Claude nutzen für:

- Komplexe Code-Analysen
- Wissenschaftliche Recherchen
- Detaillierte Dokumente erstellen
- Strategische Beratung
- Multi-Step-Reasoning

1.3 Architektur

Das hybride System verwendet einen **Smart Router**, der eingehende Anfragen intelligent zwischen Ollama und Claude verteilt:

Routing-Logik

Einfache Anfragen (kurz, faktisch, schnell) → Ollama

Komplexe Anfragen (lang, analytisch, kreativ) → Claude

Sensible Daten → Immer Ollama

Kapitel 2

Voraussetzungen

2.1 Hardware

Komponente	Anforderung
Betriebssystem	macOS 12.0 oder höher
Prozessor	Apple Silicon (M1/M2/M3) oder Intel
RAM	Mindestens 8 GB (16 GB empfohlen)
Freier Speicher	10–50 GB (je nach Modell)

Tabelle 2.1: Hardware-Anforderungen

Warnung

RAM-Empfehlung: Für optimale Performance wird mindestens 16 GB RAM empfohlen. Mit 8 GB solltest du kleinere Modelle (7B) nutzen.

2.2 Software

Erforderlich:

Homebrew installiert

Node.js v18+ installiert

Clawdbot bereits eingerichtet

2.2.1 Voraussetzungen prüfen

```
1 # Homebrew installiert?  
2 brew --version  
3  
4 # Node.js installiert?  
5 node --version  
6  
7 # Clawdbot l u f t?
```

```
8 ps aux | grep clawdbot
```

Listing 2.1: Prüf-Befehle

Kapitel 3

Installation

3.1 Ollama installieren

3.1.1 Via Homebrew (empfohlen)

```
1 brew install ollama
```

Listing 3.1: Ollama Installation

3.1.2 Alternativer Download

1. Besuche: <https://ollama.com/download>
2. Lade **Ollama für macOS** herunter
3. Öffne die **.dmg**-Datei
4. Ziehe Ollama in den Programme-Ordner

3.2 Ollama Service starten

```
1 # Im Hintergrund starten
2 ollama serve &
3
4 # Oder als LaunchAgent (automatischer Start)
5 brew services start ollama
```

Listing 3.2: Service starten

Verifizierung

Prüfe, ob Ollama läuft:

```
1 curl http://localhost:11434/api/version
```

Erwartete Ausgabe: {"version":"0.15.2"}

3.3 Qwen Modell herunterladen

3.3.1 Modell-Auswahl

Modell	Größe	RAM	Speed	Qualität	Anwendung
qwen2.5:3b	2 GB	8 GB			WhatsApp, Einfach
qwen2.5:7b	4.7 GB	16 GB			Allgemein
qwen2.5:14b	8.9 GB	32 GB			Standard
qwen2.5:32b	20 GB	64 GB			Beste Qualität

Tabelle 3.1: Qwen Modell-Varianten

3.3.2 Download-Befehle

```

1 # Qwen 2.5 14B (empfohlen)
2 ollama pull qwen2.5:14b
3
4 # Kleinere Varianten
5 ollama pull qwen2.5:7b      # Für 16 GB RAM
6 ollama pull qwen2.5:3b     # Für 8 GB RAM
7
8 # Größere Variante
9 ollama pull qwen2.5:32b    # Für 64 GB RAM

```

Listing 3.3: Modell herunterladen

3.4 Installation verifizieren

```

1 ollama list

```

Listing 3.4: Installierte Modelle prüfen

Erwartete Ausgabe:

NAME	ID	SIZE	MODIFIED
qwen2.5:14b	abc123def456	8.9 GB	2 minutes ago

Kapitel 4

Konfiguration

4.1 Clawdbot konfigurieren

Öffne die Clawdbot-Konfiguration:

```
1 nano /Users/mac/.clawdbot/clawdbot.json
```

4.2 Hybride LLM-Konfiguration

Füge folgende Konfiguration ein:

```
1 {
2   "llm": {
3     "provider": "hybrid",
4     "hybrid": {
5       "router": "smart",
6       "providers": {
7         "ollama": {
8           "baseUrl": "http://localhost:11434",
9           "model": "qwen2.5:14b",
10          "enabled": true,
11          "priority": 1,
12          "useFor": ["simple", "translation",
13                   "summarize", "whatsapp"]
14        },
15        "claude": {
16          "apiKey": "DEIN_CLAUDE_API_KEY",
17          "model": "claude-opus-4-20250514",
18          "enabled": true,
19          "priority": 2,
20          "useFor": ["complex", "code",
21                   "analysis", "research"]
22        }
23      },
24      "routing": {
25        "simple": "ollama",
26        "complex": "claude",
27        "translation": "ollama",
28        "code": "claude",
29        "summarize": "ollama",
30        "analysis": "claude",
```

```
31     "research": "claude",  
32     "whatsapp": "ollama",  
33     "default": "claude"  
34 }  
35 }  
36 }  
37 }
```

Listing 4.1: Hybrid-LLM-Konfiguration

4.3 Gateway neu starten

```
1 launchctl kickstart -k gui/501/com.clawdbot.gateway
```

Kapitel 5

Qwen Plugin

5.1 Plugin erstellen

```
1 mkdir -p /Users/mac/.clawdbot/skills/qwen
```

Listing 5.1: Qwen Skill-Verzeichnis erstellen

5.2 SKILL.md erstellen

```
1 cat > /Users/mac/.clawdbot/skills/qwen/SKILL.md << 'SKILL'
2 ---
3 name: qwen
4 description: Lokales Qwen 2.5 LLM via Ollama
5 license: MIT
6 ---
7
8 # Qwen Ollama Integration
9
10 Nutzt das lokal laufende Qwen 2.5 Modell.
11
12 ## Funktionen
13 - Offline-Betrieb
14 - Kostenlos
15 - Schnell
16 - Privacy-freundlich
17
18 ## Befehle
19 /qwen [frage] - Nutze Qwen
20 /qwen-info - Zeige Status
21 SKILL
```

5.3 package.json erstellen

```
1 {
2   "name": "@clawdbot-skills/qwen",
3   "version": "1.0.0",
4   "description": "Qwen 2.5 Ollama Integration",
```

```
5  "clawdbot": {  
6    "type": "skill",  
7    "category": "llm",  
8    "tags": ["llm", "ollama", "qwen"],  
9    "author": "Clawdbot"  
10 }  
11 }
```

Listing 5.2: Skill Metadaten

5.4 Skill aktivieren

In `clawdbot.json` hinzufügen:

```
1 {  
2   "skills": {  
3     "enabled": true,  
4     "paths": ["/Users/mac/.clawdbot/skills"],  
5     "active": ["qwen"]  
6   }  
7 }
```

Kapitel 6

Testen

6.1 Ollama direkt testen

```
1 curl http://localhost:11434/api/generate -d '{
2   "model": "qwen2.5:14b",
3   "prompt": "Was ist 2+2?",
4   "stream": false
5 }'
```

Listing 6.1: Test mit curl

Erwartete Antwort:

```
{
  "response": "2+2 ergibt 4.",
  "done": true
}
```

6.2 Interaktiver Test

```
1 ollama run qwen2.5:14b "Was ist 5+5?"
```

6.3 Clawdbot-Integration testen

6.3.1 Via WhatsApp

Einfache Anfrage (Ollama)

Du: Hallo, was ist 10-3?

Bot: (via Qwen) 10-3 ergibt 7.

Komplexe Anfrage (Claude)

Du: Schreibe eine detaillierte Analyse der Blockchain-Technologie

Bot: (via Claude) [Ausführliche Antwort...]

Kapitel 7

Troubleshooting

7.1 Connection refused

Problem

Fehler: Connection refused beim Zugriff auf Ollama

Ursache: Ollama läuft nicht

Lösung:

```
1 # Ollama starten
2 ollama serve &
3
4 # Oder als Service
5 brew services start ollama
6
7 # Status pr fen
8 curl http://localhost:11434/api/version
```

7.2 Model not found

Problem

Fehler: Model 'qwen2.5:14b' not found

Ursache: Modell nicht heruntergeladen

Lösung:

```
1 # Modell herunterladen
2 ollama pull qwen2.5:14b
3
4 # Installierte Modelle pr fen
5 ollama list
```

7.3 Ollama ist langsam

Performance-Probleme

Ollama antwortet sehr langsam oder stockt

Lösungen:

1. Kleineres Modell nutzen

```
1 ollama pull qwen2.5:7b
```

2. Andere Apps schließen

```
1 # RAM-Nutzung prüfen  
2 top -o mem
```

3. Intel Mac: Nutze 3B oder 7B Modell

7.4 Routing funktioniert nicht

Problem

Clawdbot nutzt immer Claude statt Ollama

Lösung: Routing-Konfiguration prüfen

```
1 {  
2   "routing": {  
3     "whatsapp": "ollama", // Muss auf "ollama" stehen  
4     "default": "claude"  
5   }  
6 }
```

Kapitel 8

Erweiterte Nutzung

8.1 Weitere Modelle

8.1.1 Llama 3.1 (Meta)

```
1 ollama pull llama3.1:8b
```

8.1.2 Mistral (MistralAI)

```
1 ollama pull mistral:7b
```

8.1.3 CodeLlama (Code-spezialisiert)

```
1 ollama pull codellama:13b
```

8.2 Multi-Modell-Setup

Verschiedene Modelle für verschiedene Aufgaben:

```
1 {
2   "providers": {
3     "ollama-qwen": {
4       "model": "qwen2.5:14b",
5       "useFor": ["whatsapp", "translation"]
6     },
7     "ollama-codellama": {
8       "model": "codellama:13b",
9       "useFor": ["code", "programming"]
10    },
11    "ollama-llama": {
12      "model": "llama3.1:8b",
13      "useFor": ["simple", "summarize"]
14    },
15    "claude": {
16      "model": "claude-opus-4-20250514",
17      "useFor": ["complex", "research"]
18    }
19  }
```

```

18     }
19   }
20 }

```

Listing 8.1: Multi-Modell-Konfiguration

8.3 Performance-Optimierung

8.3.1 Ollama-Konfiguration

```

1 mkdir -p ~/.ollama
2 cat > ~/.ollama/config.json << 'CONFIG'
3 {
4   "num_parallel": 4,
5   "num_ctx": 8192,
6   "num_gpu": 1,
7   "num_thread": 8
8 }
9 CONFIG

```

Listing 8.2: Ollama Config erstellen

8.3.2 Parameter-Erklärung

Parameter	Bereich	Beschreibung
num_parallel	1–8	Parallele Anfragen
num_ctx	2048–32768	Kontext-Größe (Tokens)
num_thread	1–16	CPU-Threads (= Anzahl Kerne)

Tabelle 8.1: Ollama Performance-Parameter

Kapitel 9

Sicherheit & Privacy

9.1 Warum Ollama sicherer ist

Privacy-Vorteile

- **Keine Cloud:** Daten verlassen niemals deinen Mac
- **Kein Tracking:** Ollama loggt keine Anfragen
- **Open Source:** Code ist öffentlich einsehbar
- **Offline:** Funktioniert ohne Internet

9.2 Best Practices

Sensible Daten: Immer Ollama statt Claude nutzen

Firmen-Daten: Routing auf `ollama` setzen

Passwörter: Niemals in LLM-Prompts eingeben

Logs: Regelmäßig löschen

9.3 Logs löschen

```
1 # Ollama-Logs löschen
2 rm -rf ~/.ollama/logs/*
3
4 # Clawdbot-Logs löschen
5 rm -rf /Users/mac/.clawdbot/logs/*.log
```


Anhang A

Befehls-Referenz

A.1 Ollama-Befehle

Befehl	Beschreibung
<code>brew services start ollama</code>	Service starten
<code>brew services stop ollama</code>	Service stoppen
<code>brew services list</code>	Service-Status
<code>ollama list</code>	Modelle auflisten
<code>ollama pull <model></code>	Modell herunterladen
<code>ollama run <model> "text"</code>	Modell testen
<code>ollama rm <model></code>	Modell löschen
<code>curl localhost:11434/api/version</code>	API-Version

Tabelle A.1: Ollama Befehls-Übersicht

A.2 Clawdbot-Befehle

Befehl	Beschreibung
<code>launchctl kickstart -k ...</code>	Gateway neu starten
<code>launchctl stop ...</code>	Gateway stoppen
<code>tail -f logs/gateway.log</code>	Logs anzeigen
<code>cat clawdbot.json grep llm</code>	Config prüfen
<code>ls -la skills/</code>	Skills auflisten

Tabelle A.2: Clawdbot Befehls-Übersicht

Anhang B

Installations-Checklist

B.1 Vor der Installation

- macOS 12.0+
- Homebrew installiert
- Node.js 18+ installiert
- Clawdbot läuft
- Mindestens 8 GB RAM
- 10–50 GB freier Speicher

B.2 Installation

- Ollama installiert
- Ollama Service läuft
- Qwen-Modell heruntergeladen
- Hybrid-Config eingefügt
- Qwen Skill erstellt
- Gateway neu gestartet

B.3 Tests

- Ollama API erreichbar
- Qwen-Modell läuft
- Clawdbot nutzt Ollama
- Hybrides Routing funktioniert
- Logs zeigen keine Fehler