# The AI Agent Builder's Playbook

## From Zero to Production AI Agents

*By someone who actually built them.*

---

# Introduction

You've seen the demos. AI agents booking flights, writing code, making phone calls. Most of what you've seen is vaporware. But real, working AI agents? They exist. I've built them.

This guide is the distillation of building: - A voice AI that makes real phone calls and books meetings - A crypto trading bot running 24/7 - Production LLM applications handling real users

No theory. No hype. Just what works.

---

# Chapter 1: What AI Agents Actually Are

## The Hype vs Reality

**The hype:** Autonomous systems that think for themselves and complete complex tasks without human intervention.

**The reality:** Sophisticated loops that: 1. Take input 2. Decide what to do (using an LLM) 3. Execute actions (using tools) 4. Observe results 5. Repeat until done

That's it. The magic is in the implementation details.

## The Core Loop

```
while not done:
    # 1. Get current state
    observation = get_state()

    # 2. Ask LLM what to do
    action = llm.decide(observation, tools, history)

    # 3. Execute the action
    result = execute(action)

    # 4. Check if we're done
    done = check_completion(result)
```

Every agent, from AutoGPT to Claude Computer Use, is a variation of this loop.

## What Makes Agents Hard

1. **Reliability** - LLMs hallucinate. Tools fail. Networks timeout.
2. **State management** - Keeping track of what happened across many steps.
3. **Cost** - Each LLM call costs money. Loops can get expensive.
4. **Latency** - Users don't want to wait 30 seconds for a response.

This guide teaches you to handle all of these.

---

# Chapter 2: The Modern Agent Stack

## The Tools That Matter

### LangGraph

The best framework for building agents in 2025-2026. Why: - Explicit state machines (you control the flow) - Built-in persistence (state survives restarts) - Human-in-the-loop patterns - Production-ready with LangGraph Platform

```
from langgraph.graph import StateGraph

# Define your agent as a graph
graph = StateGraph(AgentState)
graph.add_node("think", thinking_node)
graph.add_node("act", action_node)
graph.add_edge("think", "act")
graph.add_conditional_edges("act", should_continue)
```

## Claude (Anthropic)

Best for agents because: - Follows instructions precisely - Less likely to hallucinate tool calls - Understands complex multi-step tasks - Has built-in tool use

## The Tool Pattern

```
tools = [
    {
        "name": "search_web",
        "description": "Search the internet for current information",
        "parameters": {
            "query": {"type": "string", "description": "Search query"}
        }
    },
    {
        "name": "send_email",
        "description": "Send an email to someone",
        "parameters": {
            "to": {"type": "string"},
            "subject": {"type": "string"},
            "body": {"type": "string"}
        }
    }
]
```

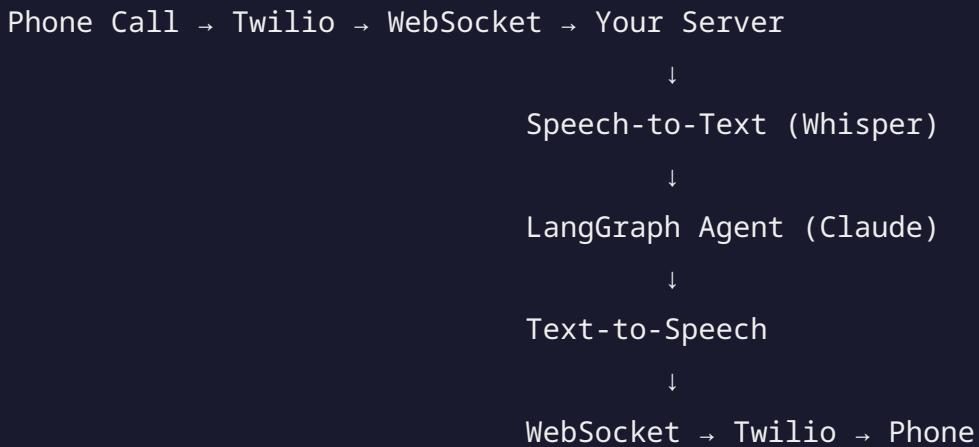The LLM sees these tool definitions and decides when to use them.

# Chapter 3: Building a Voice Agent

This is a real system I built. It: - Makes outbound phone calls via Twilio - Has natural conversations using speech-to-text and text-to-speech - Books meetings on Google Calendar - Sends SMS confirmations

## Architecture

```
Phone Call → Twilio → WebSocket → Your Server
                          ↓
                  Speech-to-Text (Whisper)
                          ↓
                  LangGraph Agent (Claude)
                          ↓
                  Text-to-Speech
                          ↓
              WebSocket → Twilio → Phone
```

## The Key Insight: Turn-Taking

Voice agents fail when they talk over people. The solution:

1. **Silence detection** - Know when the human stopped talking
2. **Strict response rules** - Ask ONE question, then STOP
3. **Interruption handling** - If they start talking, shut up immediately

```
# In your agent prompt:
"""
CRITICAL RULES:
1. Ask only ONE question per response
2. STOP talking immediately after the question mark
```

```
    3. If user interrupts, stop generating and listen
    """
```

## Real Code: The Audio Pipeline

```python
async def handle_audio_stream(websocket):
    audio_buffer = []

    async for message in websocket:
        # Accumulate audio chunks
        audio_buffer.append(message)

        # Detect silence (VAD)
        if is_silence(message, threshold=0.5):
            # User stopped talking - process
            transcript = whisper.transcribe(audio_buffer)

            # Get agent response
            response = await agent.process(transcript)

            # Convert to speech and send back
            audio = tts.synthesize(response)
            await websocket.send(audio)

            audio_buffer = []
```

# Chapter 4: Building a Trading Bot

Trading bots are agents with financial stakes. Every mistake costs money.

# The Strategy: EMA + PSAR

A profitable strategy I've run: - **EMA crossover** for trend direction - **Parabolic SAR** for entry/exit timing - **ATR** for position sizing

```python
def get_signal(df):
    # Calculate indicators
    df['ema_fast'] = ta.ema(df['close'], 5)
    df['ema_slow'] = ta.ema(df['close'], 20)
    df['psar'] = ta.psar(df['high'], df['low'], df['close'])

    # Buy signal: fast EMA crosses above slow, PSAR below price
    buy = (df['ema_fast'] > df['ema_slow']) and (df['psar'] < df['close'])

    # Sell signal: opposite
    sell = (df['ema_fast'] < df['ema_slow']) and (df['psar'] > df['close'])

    return 'BUY' if buy else 'SELL' if sell else None
```

# Risk Management

The difference between profit and ruin:

1. **Position sizing** - Never risk more than 1-2% per trade
2. **Stop losses** - Always have an exit plan
3. **Take profits** - Scale out at multiple levels

```python
def calculate_position_size(account_balance, risk_percent, stop_loss_pips):
    risk_amount = account_balance * (risk_percent / 100)
    position_size = risk_amount / stop_loss_pips
    return position_size
```

## The Execution Loop

```python
while market_is_open:
    # Get fresh data
    df = fetch_candles(symbol, timeframe)

    # Check for signals
    signal = get_signal(df)

    # Execute if we have a signal and no open position
    if signal and not has_open_position():
        if signal == 'BUY':
            open_long(calculate_position_size())
        else:
            open_short(calculate_position_size())

    # Manage existing positions
    manage_trailing_stops()
    check_take_profit_levels()

    # Wait for next candle
    sleep(timeframe_seconds)
```

# Chapter 5: Production Deployment

## Docker Everything

```dockerfile
FROM python:3.11-slim

WORKDIR /app
```

```
COPY requirements.txt .
RUN pip install -r requirements.txt

COPY . .

CMD ["python", "-m", "uvicorn", "main:app", "--host", "0.0.0.0"]
```

## Health Checks and Monitoring

```
@app.get("/health")
async def health():
    return {
        "status": "healthy",
        "llm_connected": await check_llm(),
        "db_connected": await check_db(),
        "last_activity": get_last_activity_time()
    }
```

## Handling Failures

```
from tenacity import retry, stop_after_attempt, wait_exponential

@retry(
    stop=stop_after_attempt(3),
    wait=wait_exponential(multiplier=1, min=4, max=10)
)
async def call_llm(prompt):
    return await llm.complete(prompt)
```

# Chapter 6: Common Mistakes

## 1. No State Persistence

Your agent crashes. All context is lost. User has to start over.

**Fix:** Use LangGraph with PostgreSQL backend.

## 2. Infinite Loops

Agent gets stuck calling the same tool repeatedly.

**Fix:** Add a step counter and max iterations.

```
if state.steps > MAX_STEPS:
    return "I've been trying for a while. Let me summarize what I found..."
```

## 3. Exposing API Keys

Hardcoding keys in code that gets committed.

**Fix:** Environment variables, always.

```
import os
api_key = os.environ['OPENAI_API_KEY']  # Never hardcode
```

## 4. Not Handling Rate Limits

**Fix:** Exponential backoff + request queuing.

## 5. Trusting LLM Output

LLMs lie. They confidently make up tool parameters.

**Fix:** Validate everything.

```python
def execute_tool(tool_call):
    # Validate parameters before execution
    if tool_call.name == "send_email":
        if not is_valid_email(tool_call.params['to']):
            return "Invalid email address"

    # Now execute
    return tools[tool_call.name](**tool_call.params)
```

# Conclusion

AI agents are not magic. They're engineering.

The companies winning with AI agents aren't using secret techniques. They're:
1. Building robust loops 2. Managing state properly 3. Handling failures gracefully 4. Iterating based on real usage

Now you have the playbook. Build something.

# Appendix: Resources

- LangGraph docs: https://langchain-ai.github.io/langgraph/
- Anthropic Claude docs: https://docs.anthropic.com/
- Twilio Voice: https://www.twilio.com/docs/voice

- MetaTrader5 Python: https://www.mql5.com/en/docs/python_metatrader5

---

*This guide was written by an AI engineer who ships production agents.*