# API & Queue Management Service

<div style="text-align: right;">`</> Python`</div>

```python
# api_service.py
from fastapi import FastAPI, WebSocket, Depends
from sqlalchemy.orm import Session
import asyncio
import httpx
import logging

class APIService:
    def __init__(self):
        self.app = FastAPI(title="Baker Group LLM API")
        self.inference_client =
InferenceClient("http://localhost:8001")
        self.queue_manager =
SQLiteQueueManager("queue.db")
        self.websocket_manager = WebSocketManager()

    @self.app.post("/api/v1/submit-request")
    async def submit_request(self, request:
RequestModel,
                             user: UserContext =
Depends(get_current_user)):
        """Submit request to processing queue"""
        try:
            # Validate request and check user limits
            await self._validate_request(request, user)

            # Queue request with priority
            request_id = await
self.queue_manager.enqueue_request(
                user_context=user,
                request_data=request,
                priority=PRIORITY_LEVELS[user.role]
            )

            # Notify inference service of new request
            await
self.inference_client.notify_new_request()

            # Send real-time update to user
            await
self.websocket_manager.send_user_update(
                user.id, {"request_id": request_id,
"status": "queued"}
```

```python
36                )
37
38                return {"request_id": request_id,
     "estimated_wait": await
     self._estimate_wait_time(user.role)}
39
40            except Exception as e:
41                logger.error(f"Request submission failed:
     {e}")
42                raise HTTPException(status_code=500,
     detail="Request processing failed")
43
44        @self.app.get("/api/v1/request-
     status/{request_id}")
45        async def get_request_status(self, request_id: str,
46                                      user: UserContext =
     Depends(get_current_user)):
47            """Get current request status and results"""
48            request_status = await
     self.queue_manager.get_request_status(request_id,
     user.id)
49
50            if request_status["status"] == "completed":
51                # Retrieve result from response storage
52                result = await
     self.queue_manager.get_response_data(request_id)
53                return {"status": "completed", "result":
     result}
54
55            return request_status
56
57        @self.app.websocket("/ws/{user_id}")
58        async def websocket_endpoint(self, websocket:
     WebSocket, user_id: str):
59            """Real-time updates for request status"""
60            await self.websocket_manager.connect(websocket,
     user_id)
61            try:
62                while True:
63                    # Keep connection alive and handle
     ping/pong
64                    await websocket.receive_text()
65            except WebSocketDisconnect:
66
     self.websocket_manager.disconnect(websocket, user_id)
67
68 class InferenceClient:
69     """Client for communicating with inference
     service"""
70     def __init__(self, inference_url: str):
71         self.client =
     httpx.AsyncClient(base_url=inference_url, timeout=300.0)
```

```python
72
73      async def notify_new_request(self):
74          """Notify inference service of new queued
    requests"""
75          try:
76              await self.client.post("/internal/process-
    queue")
77          except httpx.RequestError as e:
78              logger.warning(f"Failed to notify inference
    service: {e}")
79
80      async def health_check(self) -> bool:
81          """Check if inference service is available"""
82          try:
83              response = await self.client.get("/health")
84              return response.status_code == 200
85          except:
86              return False
87
88
89
```

## LLM Inference Service

</> Python

```python
1   # inference_service.py
2   from fastapi import FastAPI
3   import asyncio
4   import sqlite3
5   from contextlib import asynccontextmanager
6   import time
7
8   class InferenceService:
9       def __init__(self):
10          self.app = FastAPI(title="Baker Group Inference
    Service")
11          self.queue_manager =
    SQLiteQueueManager("queue.db")
12          self.model_manager = ModelResourceManager()
13          self.processing_loop_task = None
14
15      @asynccontextmanager
16      async def lifespan(self, app: FastAPI):
17          """Manage service lifecycle"""
18          # Startup
19          await self.model_manager.load_models()
20          self.processing_loop_task =
    asyncio.create_task(self._processing_loop())
21          yield
```

```
22          # Shutdown
23          if self.processing_loop_task:
24              self.processing_loop_task.cancel()
25          await self.model_manager.cleanup()
26
27      async def _processing_loop(self):
28          """Main inference processing loop"""
29          while True:
30              try:
31                  # Get next highest priority request
32                  request = await
    self.queue_manager.get_next_request()
33
34                  if request:
35                      await
    self._process_inference_request(request)
36                  else:
37                      await asyncio.sleep(1)  # Brief
    pause if no requests
38
39              except Exception as e:
40                  logger.error(f"Processing loop error:
    {e}")
41                  await asyncio.sleep(5)  # Error
    recovery delay
42
43      async def _process_inference_request(self, request:
    QueuedRequest):
44          """Process individual inference request"""
45          try:
46              # Update request status to processing
47              await
    self.queue_manager.update_request_status(request.id,
    "processing")
48
49              # Determine appropriate model and acquire
    resource lock
50              if request.model_target == "gpt-oss-20b":
51                  async with
    self.model_manager.gpt_oss_lock:
52                      result = await
    self._execute_text_inference(request)
53              else:  # llama3.2-vision-11b
54                  async with
    self.model_manager.vision_lock:
55                      result = await
    self._execute_vision_inference(request)
56
57              # Store result and update status
58              await
    self.queue_manager.store_response(request.id, result)
59              await
```

```
                self.queue_manager.update_request_status(request.id,
        "completed")
60
61                    # Notify API service of completion
        (optional webhook)
62                    await self._notify_completion(request.id,
        request.user_id)
63
64            except Exception as e:
65                    await
        self.queue_manager.update_request_status(request.id,
        "failed", str(e))
66                    logger.error(f"Inference request
        {request.id} failed: {e}")
67
68        @self.app.post("/internal/process-queue")
69        async def trigger_queue_processing(self):
70            """Internal endpoint to trigger queue
        processing"""
71            # This endpoint allows API service to notify of
        new requests
72            return {"status": "processing triggered"}
73
74        @self.app.get("/health")
75        async def health_check(self):
76            """Health check endpoint"""
77            model_status = await
        self.model_manager.get_model_status()
78            gpu_status = await
        self.model_manager.get_gpu_status()
79
80            return {
81                "status": "healthy" if
        model_status["loaded"] else "degraded",
82                "models": model_status,
83                "gpu": gpu_status,
84                "queue_depth": await
        self.queue_manager.get_queue_depth()
85            }
86
87 class ModelResourceManager:
88    """Manages Ollama models and GPU resources"""
89    def __init__(self):
90        self.gpt_oss_lock = asyncio.Semaphore(1)
91        self.vision_lock = asyncio.Semaphore(1)
92        self.ollama_client = OllamaClient()
93
94    async def load_models(self):
95        """Initialize and load both models"""
96        await self.ollama_client.pull_model("gpt-
        oss:20b")
97        await self.ollama_client.pull_model("llama3.2-
```

```
       vision:11b")
 98            logger.info("Models loaded successfully")
 99
100      async def get_model_status(self) -> dict:
101          """Get current model loading status"""
102          return {
103              "loaded": True,  # Check actual model
      status
104              "gpt_oss_memory": "14GB",
105              "vision_memory": "12GB",
106              "available_vram": "6GB"
107          }
108
109
110
```

inter-service API

```
                                              </> Python

 1   # Communication between API service and Inference
     service
 2   class ServiceCommunication:
 3       def __init__(self):
 4           self.api_service_url = "http://localhost:8000"
 5           self.inference_service_url =
     "http://localhost:8001"
 6
 7       # API Service -> Inference Service
 8       async def notify_new_request(self):
 9           """API notifies inference of new queued
     requests"""
10           endpoint = f"
     {self.inference_service_url}/internal/process-queue"
11
12       async def check_inference_health(self):
13           """API checks if inference service is
     available"""
14           endpoint = f"
     {self.inference_service_url}/health"
15
16       # Inference Service -> API Service (optional
     webhook)
17       async def notify_completion(self, request_id: str,
     user_id: str):
18           """Inference notifies API when request
     completes"""
19           endpoint = f"
     {self.api_service_url}/internal/request-completed"
20           payload = {"request_id": request_id, "user_id":
     user_id}
21
```

```
22
23
```

healthcheck service

</> Python

```python
class ServiceHealthMonitor:
    def __init__(self):
        self.api_service = "http://localhost:8000"
        self.inference_service = "http://localhost:8001"

    async def monitor_services(self):
        """Continuous service health monitoring"""
        while True:
            api_health = await self._check_service_health(self.api_service)
            inference_health = await self._check_service_health(self.inference_service)

            if not api_health:
                await self._alert_service_down("API Service")
            if not inference_health:
                await self._alert_service_down("Inference Service")

            await asyncio.sleep(30)  # Check every 30 seconds

    async def get_system_status(self):
        """Get overall system health status"""
        return {
            "api_service": await self._detailed_health_check(self.api_service),
            "inference_service": await self._detailed_health_check(self.inference_service),
            "queue_metrics": await self._get_queue_metrics(),
            "gpu_status": await self._get_gpu_metrics()
        }
```