# Resume — high-level critique (short and honest)

- **Why backend/production focus vs frontend or ML?**
  Because you want systems responsibility: reliability, scale, infra, ops, and glue-code that actually runs in production. Good reason. Pushback: if you care about user-facing polish or ML research novelty, saying "I focus on backend" must come with evidence you shipped systems, not just prototypes.

- **Sentence an interviewer will pick apart first:**
  Anything that says **"production-ready"** or a precise percentile latency (e.g., "2.2s median warm-request latency") — those are measurable and invite follow-ups.

- **If forced to prove one claim live:**
  Ask to reproduce the **cache warmup that yielded a 60% improvement** or show logs/benchmarks for the 2.2s median warm-request latency. Why? They are empirical, reproducible, and hard to fake on the spot without evidence.

- **Summary / high-level claims you must defend:**

  - "Production-ready" means you handled CI/CD, secrets, monitoring, graceful degradation, and recovery.
  - "Applied Generative AI" must mean implemented RAG or model calls in user-facing workflows, not experiments in a notebook.

# Definitions, short and precise

- **Production-ready (one precise sentence):**
  A system is production-ready when it consistently meets defined SLAs under realistic load, has automated CI/CD, monitoring and alerting, secure secret management, and documented rollback and recovery procedures.
- **Applied Generative AI (one sentence):**
  Using LLMs or generative models as integrated, production-facing components that ingest user or domain data, perform retrieval/augmentation or generation, and return responses with measurable latency, cost, and confidence controls.

# Which parts of production matter most — and which do you actually implement?

- **Most important:** correctness, availability, latency percentiles (p95/p99), observability, security, and cost control.
- **Common minimal implementable subset for projects:** CI/CD, automated tests, basic monitoring/logging/tracing, secrets in vault or env, simple autoscaling and health checks, graceful degradation (rate limit or fallback). That's what you should actually build if you claim "production-ready."

# Tradeoffs you accepted (prod vs research)

- Favored deterministic and debuggable over pure performance optimizations.
- Limited model calls for cost and latency instead of maximal contextualization.
- Chose simpler, maintainable infra to meet SLA rather than bleeding-edge micro-optimizations. Correct tradeoff for internships and startups.

# Education and fundamentals

- **DSA recently used in a project:** implement an LRU cache (doubly linked list + hashmap) to speed repeated queries in Chatify presence/messaging; complexity O(1) per op. If asked, show the code and tests.

- **Systems concept that influenced design (PanditAI/Chatify):** backpressure and rate limiting (token bucket) influenced API gateway and worker queue sizing. Concrete effect: limit concurrent LLM calls per user and back off gracefully.

- **How networks/OS change backend design:**

  - **Networks:** latency and bandwidth shape batching, caching, and retry strategy. Design for eventual consistency across services.
  - **OS:** process vs thread vs async choices. For Node, avoid blocking CPU-bound ops. For Python FastAPI, choose uvicorn+workers or Gunicorn+uvicorn-worker based on CPU vs IO profile.

**Gaps to fill:** deeper OS kernel internals, distributed consensus (Raft), and advanced capacity planning. Say that; interviewer will like the honesty.

# PanditAI — end-to-end, practical stuff

- **Request path (Next.js → FastAPI) (concise):**
  Browser → Next.js API Route or client fetch → Auth cookie/JWT → API gateway / nginx → FastAPI (uvicorn) route → middleware (auth/validation) → controller → service (business logic + retrieval + LLM orchestration) → Neo4j / Mongo / Postgres / Redis → response → Next.js renders.

- **Sessions & auth:** JWT access token in HttpOnly cookie; refresh-token stored as secure, rotating token in DB or Redis with revocation list. Use short access TTL and silent refresh.

- **Typical FastAPI endpoints & schemas (example):**

  - `POST /api/query`
    Request: `{ "query": "string", "user_id": "uuid", "context": {}}`
    Response: `{ "id": "uuid", "answer": "string", "sources": [{"id":"", "score":0.9}], "latency_ms": 340 }`
  - `GET /api/user/:id/profile` etc.

- **Pydantic + TypeScript drift prevention:** single-source-of-truth codegen. Use Pydantic models, export OpenAPI, generate TypeScript types via `openapi-generator` or `datamodel-code-generator`, with CI job failing on drift.

- **Neo4j vs Mongo role:**

  - Neo4j: knowledge graph, relationships, recommendations, path queries.
  - Mongo: document storage for chat logs, raw documents, user blobs.

- **Example Neo4j query + complexity:**

  ```
  MATCH (u:User {id:$uid})-[:INTERESTED_IN]->(t:Topic)-[:RELATED_TO]->
  RETURN p LIMIT 10
  ```

  Complexity: depends on node degrees; with index on `:User(id)` start is O(1), traversal cost roughly proportional to matched neighborhood size. Add indexes on high-cardinality properties, and use relationship indexes/constraints where possible.

- **LLM calls per request & rate limits:** Prefer 1–3 calls per user request (retrieval -> synthesize -> optional refinement). Manage via token-bucket at gateway, per-user quotas, and batching where possible.

- **Cold-start/warm-request latencies & tooling:** Use realistic load tests (k6 or locust) and APM (Datadog/NewRelic) to measure tail latencies. Instrument cold start by measuring first request after container scale-up; warm requests measured after keepalive.

- **Cache warm strategy & 60% improvement:** preload top N items into Redis at deploy, populate caches during low-traffic windows, and implement async prefetch on common request patterns. Proven improvement by controlled A/B benchmark: baseline median latency vs warmed-cache median.

- **Redis eviction policy chosen:** `volatile-lru` for TTLed keys or `allkeys-lru` if everything is cacheable. Choose `volatile-lru` if some keys must persist and use TTLs for read-heavy cached items.

- **Shrinking GPU image by 700MB (high-level steps):**

  1. Move to multi-stage build.
  2. Replace CUDA base with `python:3.X-slim` or a smaller runtime.
  3. Remove GPU-only libs and wheels.
  4. `pip wheel` dependencies in build stage and copy only necessary wheels.
  5. Clean apt caches and remove docs.

- **Logging & debugging tail latency spikes:** centralized structured logs (JSON) into ELK/Datadog; distributed tracing (OpenTelemetry → Jaeger) to find slow spans; correlate high-latency traces with CPU, GC, thread starvation, DB slow queries.

- **Contradictory facts from knowledge graph:** rank by provenance score, recency, and trust; return multiple sources flagged with confidence and let UI surface "conflicting sources" to user. Have a tie-breaking priority (most recent, curated sources).

- **If Neo4j unavailable for 30s:** degrade gracefully: fallback to cached results in Redis or progressively less-accurate retrieval from Mongo. Return partial results with `partial=true` and log incident.

# eVakeel — legal retrieval and production safety

- **Ingest & structure legal acts:** OCR/PDF → canonicalize text → split by section/paragraph → metadata index (jurisdiction, act id, version) → embeddings + chunk store and symbolic index (TOC + citations graph).

- **Retrieval approach:** hybrid: vector embeddings for semantic match + symbolic index / BM25 for exact citations. Combine scores to avoid hallucination.
- **Pipeline for ~50ms warm answer (microsteps):** authenticate → parse query → quick keyword match in inverted index → retrieve top-K doc chunks from vector store → lightweight reranker (TF-IDF or small encoder) → construct prompt → LLM call (if needed) → format. Warm path caches top-K results for frequent queries.
- **Context window and hallucination guard:** chunking with overlap, prompt templates with citation injection, LLM output parsing for citations, and a grounding checker that verifies citations exist and returns `low-confidence` if mismatch.
- **Handling law updates:** versioned documents with an `effective_date` and a change feed. Re-ingest changed docs, update embeddings, and invalidate caches for affected keys.
- **File uploads & parsing:** dedicated file-service, use robust parsers and queue for processing with error handling; malformed docs go to quarantine and a retry pipeline with human review.
- **Node ↔ Python concurrency and protocol:** use HTTP/gRPC or message queue (RabbitMQ/Kafka) for async heavy jobs. For synchronous low-latency calls, prefer gRPC with connection pooling.

# Chatify — real-time messaging at scale

- **Message lifecycle:** client → Socket.IO emit → gateway → auth middleware → matchmaking to instance (via sticky session or Redis adapter) → message persisted async to Mongo → Redis pub/sub for fanout → recipient instances deliver → ack back to sender → message considered durable after Mongo write and ack.
- **Socket.IO scaling:** use a Redis adapter; ensure sticky sessions on load balancer to keep a client connected to same instance for session-local state or use shared session in Redis. For reconnects, use user-to-instance mapping in Redis to route clients.
- **Redis-based presence & failover:** presence stored as ephemeral key with TTL in Redis. To avoid loss on failover use Redis Cluster with replication, or persist presence stream into Redis Streams or an auxiliary durable store. On failover, revalidate presence with heartbeat from clients.
- **Rate limiting:** token-bucket implemented with Lua script in Redis for atomic ops and strict limits. Use user-level and IP-level buckets.
- **Message durability guarantee:** ensure Mongo writes are done before acknowledging delivery (or use write-ahead log pattern). If you must ack before DB write for latency, use an append-only durable queue like Kafka then write to Mongo asynchronously and use replay to ensure durability.
- **If Redis shows user online but Mongo write fails:** keep message in retry queue; notify sender of pending delivery and attempt retry. Mark message state as `pending` until durable.
- **Preserving ordering under concurrency:** use per-conversation monotonic sequence numbers and append-only logs. Shard by conversation id to avoid cross-shard ordering issues.
- **Mongo schema tradeoff:** choose document model with `messages` collection embedding small chunks or referencing based on message size; embed recent messages for quick reads and

reference older ones to avoid large documents.

- **Deduplication:** unique message IDs assigned by client or server and Mongo unique index on `(conversationId, messageId)`.
- **Group chat fanout heuristics:** avoid naïve fanout by limiting fanout for inactive users, using compact deltas for large groups, and using selective push + pull model where clients fetch missed messages on reconnect.
- **Archiving:** cold storage of old messages in S3, and TTL deletes from Mongo with pointer to archived location.

# API design, contracts, schema

- **Naming conventions:** REST: `/api/v1/resources/{id}` and use nouns, consistent verbs for actions only when needed. Prefer consistent pluralization.

- **Versioning strategy:** URI versioning + backward compatible changes only; for breaking changes, create `/v2` and a migration guide. Support feature flags and gradual rollouts.

- **Idempotency:** accept `Idempotency-Key` header for non-idempotent ops; store key→result mapping.

- **Response codes:**

    - 400 for client validation errors with structured error body,
    - 401 for auth, 403 for forbidden,
    - 404 for not found,
    - 409 for conflict,
    - 500 for server errors.

- **Docs sync:** generate OpenAPI from code and publish auto-generated docs; enforce CI to fail on spec drift.

# Databases & persistence

- **Why Mongo + Neo4j instead of single DB:** different query patterns: document store for messages, graph DB for relationship traversal and knowledge graph. Using best tool for the query is fine; cost is complexity in transactions.

- **ACID across multiple DBs:** implement sagas or two-phase commit alternatives; prefer idempotency and compensation transactions.

- **Example Postgres index you added:**

```
CREATE INDEX idx_user_email_lower ON users (lower(email));
```

Before: full table scan on login by email for 1M rows. After: select latency drop from 120ms to 2ms median.

- **Connection pooling:** use pgbouncer for Postgres, limit pool size based on max clients, and use circuit-breaker to avoid thundering herd.

- **Default isolation level:** `READ COMMITTED` for most web apps; raise to `REPEATABLE READ` for complex invariants only when needed.

- **Backups & restore:** logical dumps and WAL shipping; restore time depends on DB size — practice restoring into a staging cluster and document RTO/RPO.

# Caching & performance

- **What to cache:** heavy read results, top-K retrievals, embeddings, computed responses. TTLs: short for user-specific, longer for static domain data.
- **Invalidation:** targeted keys, versioned keys (include data version in cache key), and publish-subscribe invalidation for multi-instance setups.
- **Measure cache hit rates:** instrument metrics in Prometheus; acceptable hit ratio depends on traffic but aim > 70% for high-value caches.
- **Serialization benchmark choices:** JSON is simplest; MessagePack reduces size/CPU; use binary for ultra-low-latency systems.
- **Latencies and percentiles:** track p50/p95/p99; outliers often from cold starts, GC pauses, or external dependency spikes.

# Real-time & concurrency specifics

- **Backpressure:** limit connections, reject or queue new requests, and throttle LLM calls with queues and worker pools.
- **Avoid Node blocking:** offload CPU tasks to worker threads or separate microservices; use streaming for large payloads.
- **FastAPI worker model:** uvicorn with multiple workers behind Gunicorn is common; choose async workers for IO heavy, sync workers for CPU tasks with proper worker count.
- **Fairness:** per-user quotas and queue-based scheduling.

# Generative AI specifics

- **Model choices:** justify based on latency/cost/accuracy. E.g., small in-house model or OpenAI shorter-latency deployment for interactive features, larger models for long-form jobs.
- **Prompt engineering:** templatize prompts, keep templating centralized, and version prompts.
- **PII protection:** mask or redact PII before sending, use encryption in transit and at rest, and avoid sending sensitive user ids to third-party LLMs where possible.
- **Hallucination detection:** citation checks, grounded prompts, and post-rerank verification. Present confidence scores and attach sources.

# Docker, CI/CD, deployments

- **Dockerfile best-practices:** multi-stage builds, pinned base images, minimize layers, remove caches. Show relevant Dockerfile in interview and explain each layer.
- **Local compose vs prod:** Compose for dev convenience; real production uses Kubernetes or managed container services. Key diffs: autoscaling, secrets management, observability.
- **Migrations & zero-downtime:** run schema migrations with small, backward compatible steps and feature flags. Use blue/green or canary deployments.
- **Rollback:** keep prior image in registry, run automated rollback on health-check failure.

# Cloud & infra

- **Why Azure (if used):** integration with AD, managed Postgres, good dev credits; but be ready to justify vs AWS/GCP with cost and feature comparisons.
- **Autoscaling triggers:** use CPU, queue length, and p95 latency. Tune scale steps and cool-downs.
- **RTO/RPO:** define for business and configure replicas accordingly.

# Security & auth

- **Tokens:** JWT for stateless access with short TTL, refresh tokens opaque and stored server-side, with revocation. For high security, use rotating refresh tokens.

- **RBAC & object level perms:** combine role checks with resource-level ACLs; enforce checks in service/business layer.
- **Frontend XSS/CSRF:** sanitize inputs, use HttpOnly cookies for JWT, enable sameSite, use CSP headers. For Next.js, escape server-side rendered content.
- **Secrets:** Azure Key Vault / HashiCorp Vault or managed secret store. Rotate periodically.

# Observability & SRE

- **Minimum metrics:** request rate, error rate, p50/p95/p99 latency, CPU/mem, DB connections, queue length.
- **Tracing:** OpenTelemetry for cross-service tracing.
- **Logs:** structured JSON to ELK or Datadog, with retention policy.
- **Alerts:** set on symptom-based metrics (error surge, latency increase), with escalation playbooks.

# Testing & quality

- **Testing pyramid:** unit tests (fast), integration tests, E2E (Playwright). For websockets: write integration tests mocking network and using headless clients.
- **Load tests:** k6 or locust with scenarios (10k concurrent users), measure latency, error rates, and resource usage.
- **Avoid flaky tests:** isolate external deps with mocks, inject time dependencies.

# Data modeling & algorithms

- **Messages schema tradeoff:** embed recent messages for fast read; reference older messages. Pros: faster reads; cons: document growth and writes.
- **Graph modeling:** nodes for Users, Documents, Topics; relations: AUTHORED, CITED, RELATED; index node ids and high-selectivity properties.
- **Algorithm example:** rate limiter with $O(1)$ per request using Redis fixed-window counter or sliding window with sorted set and periodic cleanup.

# Git, collaboration & engineering practices

- **Branching:** trunk-based for rapid delivery; feature branches with PRs for larger changes.
- **Bad PR example critique:** big diff, multiple concerns, no tests, missing changelog. Ask for smaller focused PRs with clear description.
- **Commit messages:** short subject line, longer body explaining why. Use type prefixes (feat, fix, chore) if team enforces.

# Soft skills & product

- **Role in hack projects:** show concrete decisions you made: architecture tradeoffs, performance optimizations, hiring or mentoring. If you can't list specifics, don't claim leadership.
- **Handling feedback:** describe a real incident, what you learned, and a precise change you made.

# Practical interview prompts — what you should be able to do live

- Reproduce a Docker build and explain each step. (Be ready, don't bluff.)
- Show a query to find slow Mongo queries: `db.system.profile.find({ millis: { $gt: 100 } })` after enabling profiling. Then propose index like `db.messages.createIndex({ conversationId: 1, createdAt: -1 })`.
- Design a 10k-concurrent-user Chatify test: use k6 with ramping VUs, measure p99 latency, message delivery success, DB write throughput.
- Triage a knowledge-graph timeout: check tracing to find which span blocks, identify resource exhaustion or network errors, fallback to cache, re-route traffic, and escalate.

# Final blunt advice you'll actually use in interviews

- Don't claim "production-ready" without numbers, diagrams, and a single runbook link.
- Be able to pull up one metric or log and explain what it proves.
- Show one small demo that proves behaviour: the cache warm script, the Dockerfile change, or the Postgres index result.

- Practice answering "how would you debug X?" live with a calm step-by-step checklist: reproduce, isolate, instrument, mitigate, fix, postmortem.

You wanted everything — you got a condensed product-grade cheat sheet. Don't memorize it verbatim. Practice answering these like you wrote the code. If you want, I can convert any one section into a 10-minute mock interview + exact sample answers you should deliver. But you probably want to start repping these out loud, not read another wall of text.