# Learn to

# Design Cloud

# Architecture

# Let's build a Real Time chat system

- A real-time chat application allows users to send and receive messages instantly over the internet.



WhatsApp

Facebook Messenger

WeChat

Line Messenger

Telegram

Signal

slack

Microsoft Teams

Discord

# WhatsApp

- **2009:**
  - WhatsApp Inc. was founded by former Yahoo employees

- **2013:**
  - 200 million+ users

- **2024:**
  - 2 billion+ active users
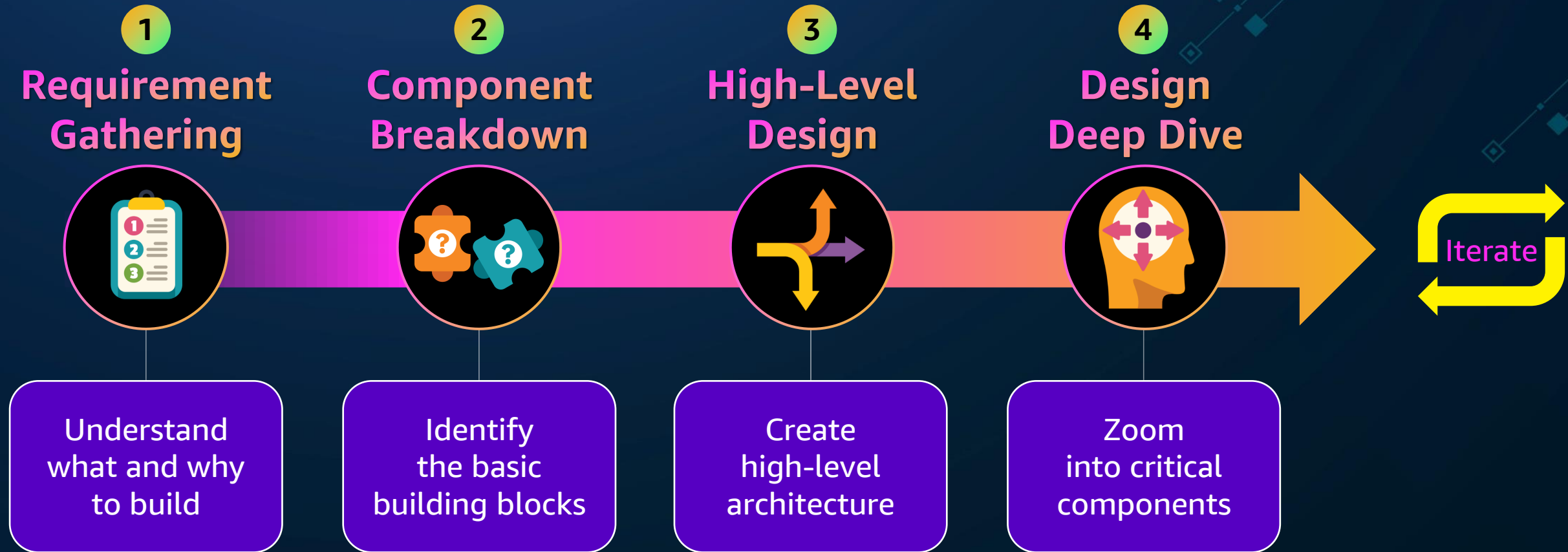


**Top features:**

- End-to-end encrypted chats
- Multi-device access support
- Free voice/video calling

# Designing a System – A simple framework

**1** Requirement Gathering

**2** Component Breakdown

**3** High-Level Design

**4** Design Deep Dive

Iterate

Understand what and why to build

Identify the basic building blocks

Create high-level architecture

Zoom into critical components

Requirement Gathering

# Functional Requirement

**1 to 1 Chat**

**Group Chat**

**Media Attachment**

**Online/ Offline status**

**Delivery Confirmation**

**Push Notifications**

# Non-functional Requirement



| **Performance** | **Scalability** | **Security** | **Reliability** | **Availability** |
|---|---|---|---|---|
| Send and receive messages in real-time | Billions of users and billions of messages | End-to-end encryption of messages | Store messages until they are delivered | Minimal downtime or interruptions |

# Out of scope

- Authentication
- Encryption
- User registration
- User profiles
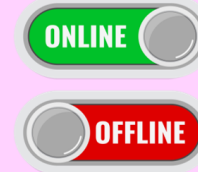- Mobile app
- Audio/Video call
- Block lists
- Multiple device support
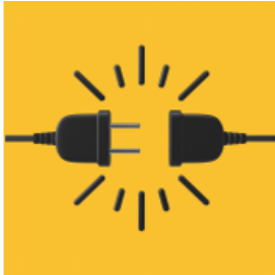
# Component Breakdown

# Core components

**Mobile App** — Front-end interface for the user

**Status Service** — Tracks user presence (online/offline)

**Connection Service** — Manages real-time connectivity

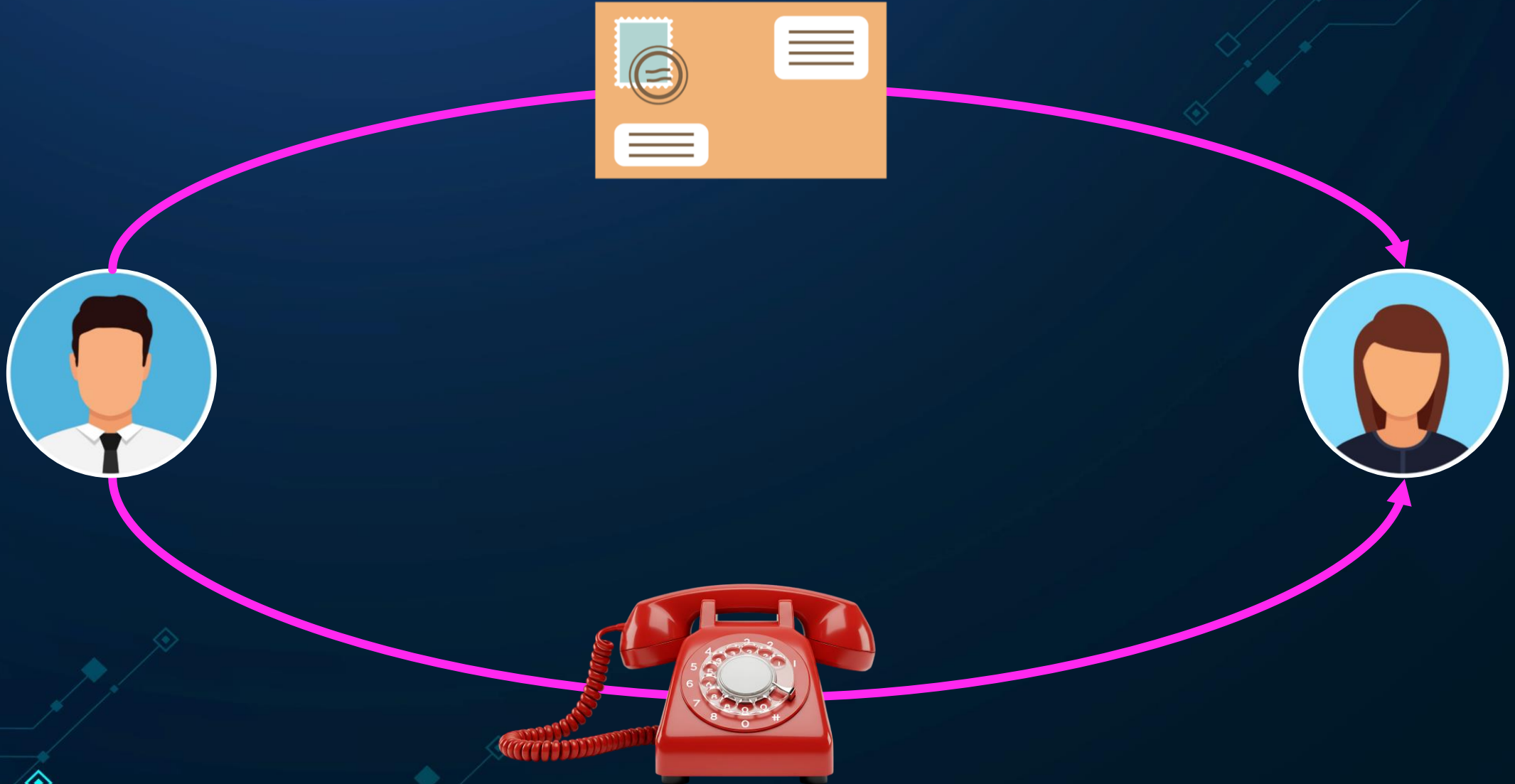**Storage Service** — Storage for chat history and media

**Message Service** — Sending and receiving of messages.

**Group Service** — Manages group chat functionality
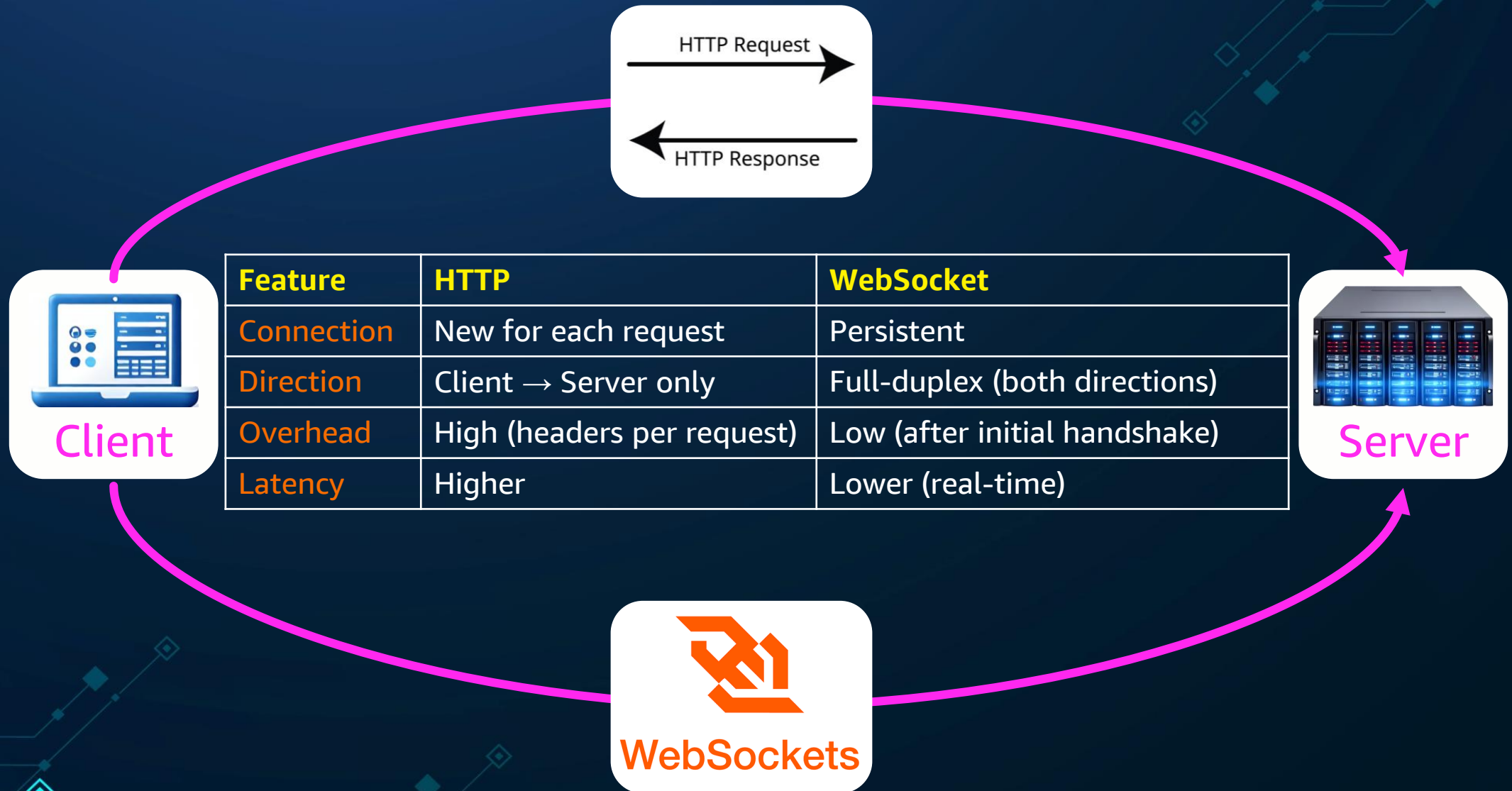
# Analogy: Sending Letters vs. A Phone Call

# What is WebSocket?

WebSocket is a communication protocol that provides full-duplex communication channels over a single, long-lived TCP connection.

# Analogy: Sending Letters vs. A Phone Call

HTTP Request →

← HTTP Response

Client

| Feature | HTTP | WebSocket |
|---|---|---|
| Connection | New for each request | Persistent |
| Direction | Client → Server only | Full-duplex (both directions) |
| Overhead | High (headers per request) | Low (after initial handshake) |
| Latency | Higher | Lower (real-time) |

Server

WebSockets

# How WebSocket Works?

| Step 1 | Handshake (via HTTP) |
| --- | --- |
| Step 2 | Persistent Connection |
| Step 3 | Transfer of Data |

Client

Server

**Common Use Cases**

| Live Chats | Online Games |
| --- | --- |
| Stock Tickers | Collaborative Editing |

# Alternative to WebSockets

- MQTT
  - MQTT stands for Message Queuing Telemetry Transport. It is a lightweight, open-source messaging protocol that is widely used in the Internet of Things (IoT) for communication between devices.

- XMPP
  - The full form of XMPP is Extensible Messaging and Presence Protocol. It is a communication protocol used for instant messaging, presence information, and contact list maintenance. XMPP is based on XML and enables near-real-time exchange of structured data between network entities.

| Feature | WebSocket | MQTT | XMPP |
|---|---|---|---|
| Communication Pattern | Full-duplex, bidirectional streams | Publish/subscribe via broker | Client-server messaging with presence |
| Data Format | Binary or text frames | Lightweight binary (topics + payload) | XML stanzas (text-based, verbose) |
| Built-in Presence | No | No | Yes, presence, contact lists, subscriptions |
| Protocol Overhead | Low (once handshake completes) | Very low, optimized for IoT/mobile | High (XML verbosity, overhead) |
| Applications (likely usage) | • Facebook Messenger (Custom WebSocket)<br>• WeChat<br>• Discord<br>• Microsoft Teams<br>• Line Messenger<br>• Signal | • Facebook Messenger (Android/iOS clients)<br>• Instagram DM | • WhatsApp - Custom XMPP "FunXMPP" |

# How to implement Web Sockets?

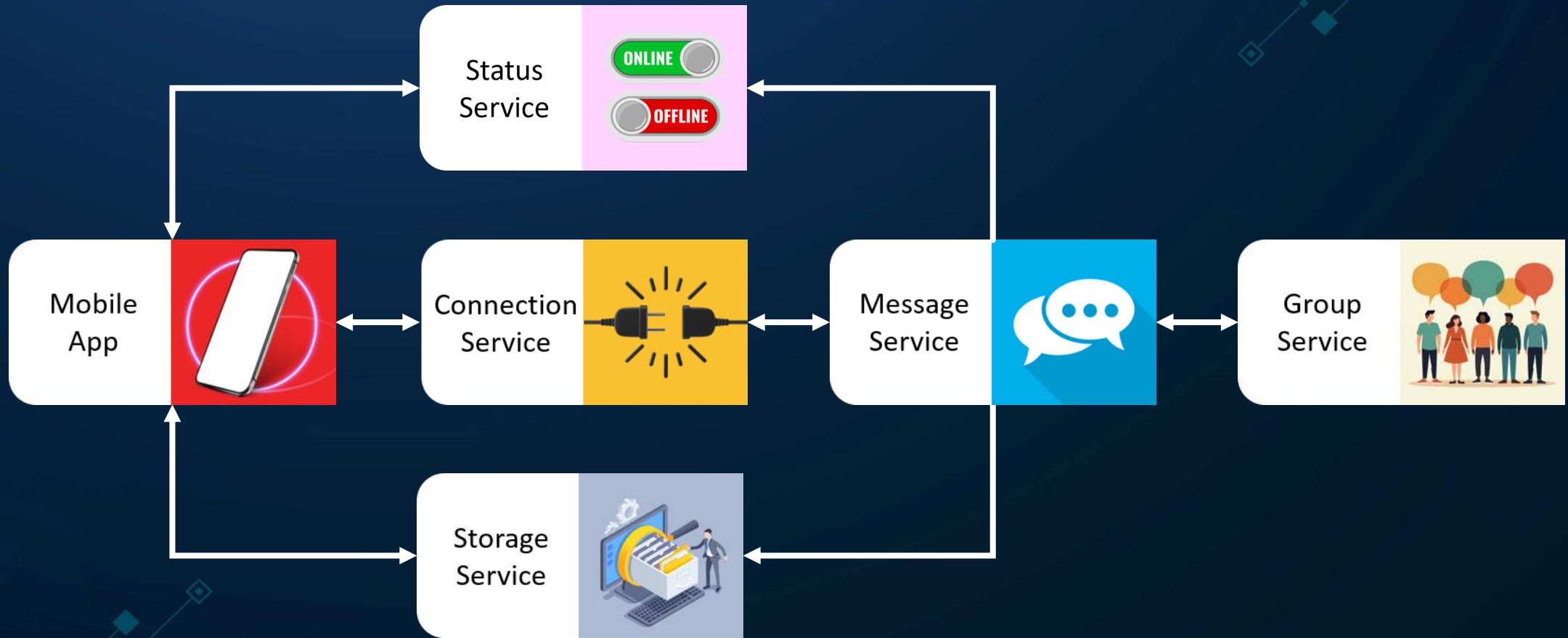| Platform | Lightweight Option<br><br>These give you direct access to the WebSocket protocol (without additional logic or structure. You'll need to manually handle reconnects, rooms, user state, etc. | High-Level Option<br><br>These libraries abstract away the raw WebSocket protocol and give you extra features out of the box, such as: Automatic reconnection Presence detection etc. |
|---|---|---|
| Java | Jetty, Tyrus | Spring WebSocket (STOMP) |
| .NET | System.Net.WebSockets | SignalR |
| Python | websockets, FastAPI | Django Channels, Socket.IO |
| Node.js | ws | Socket.IO |

High-Level Design

# High-level design
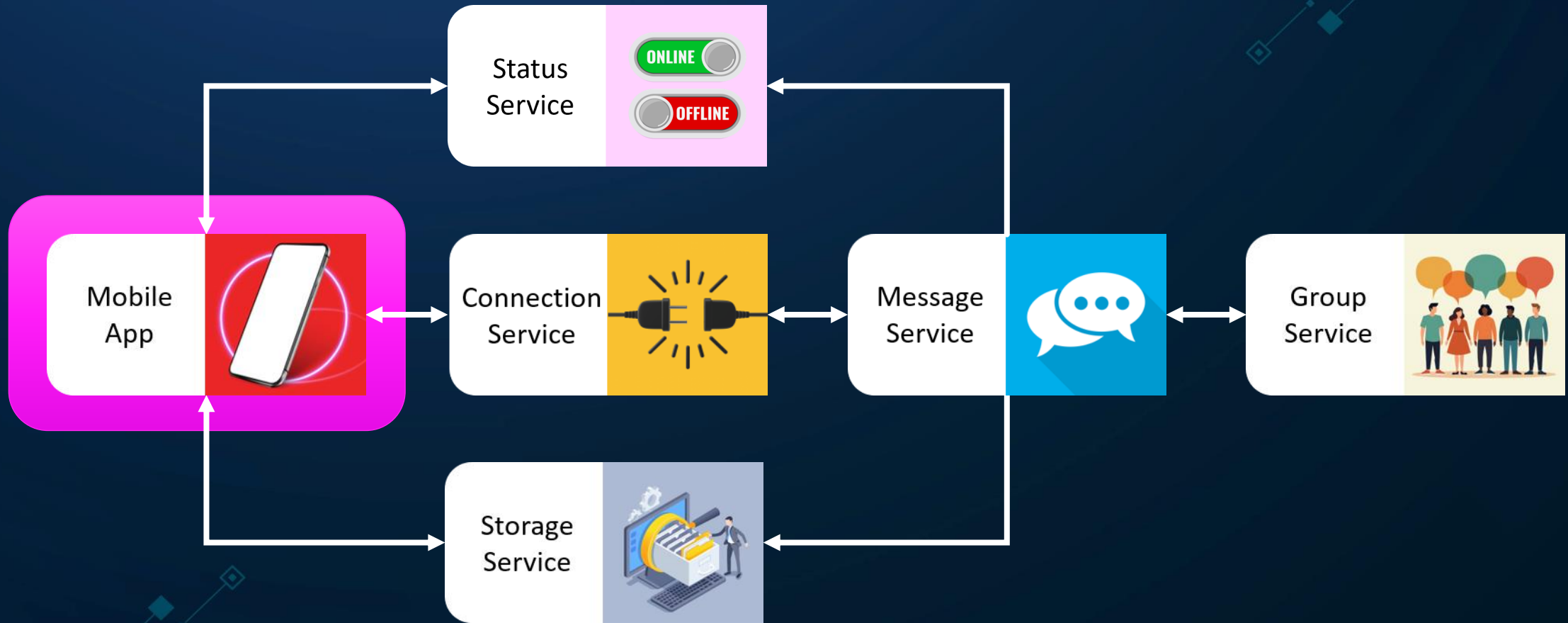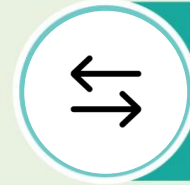
# Mobile App

# Mobile App
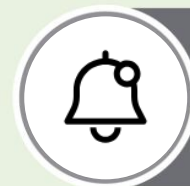


- Authentication & User Management
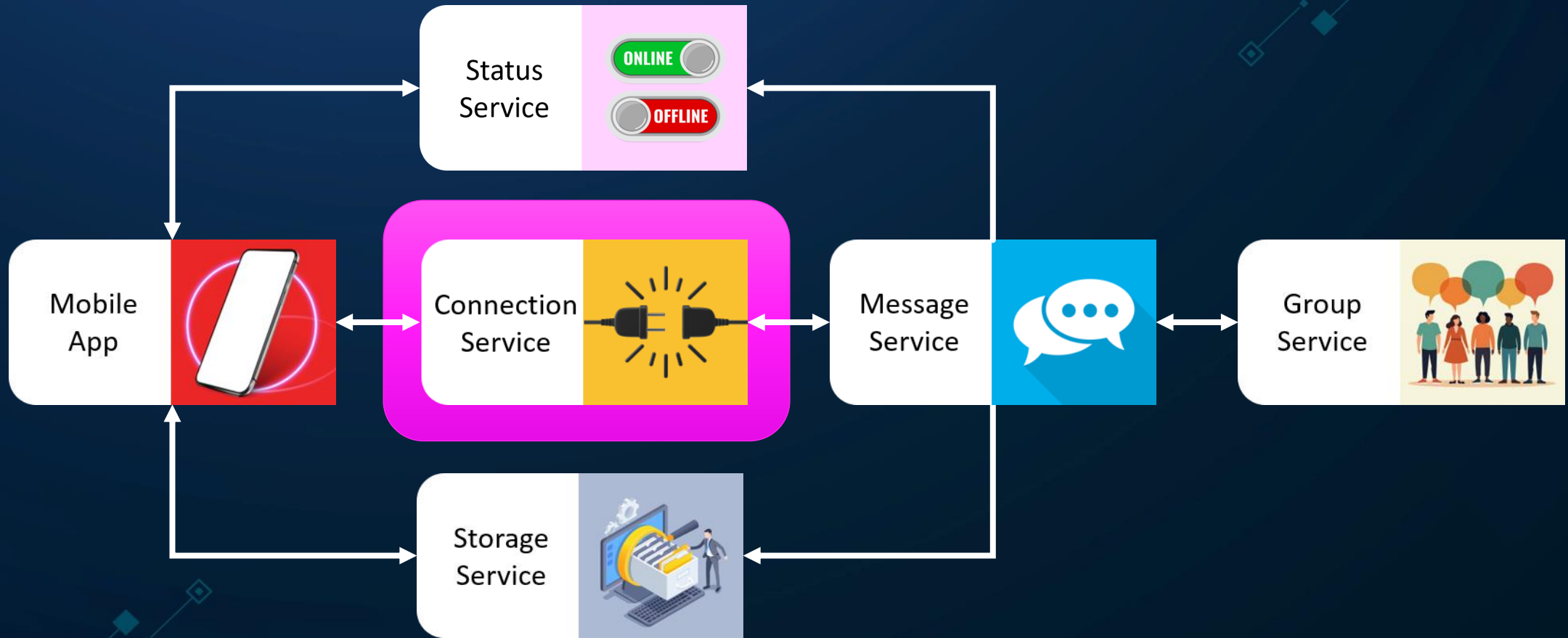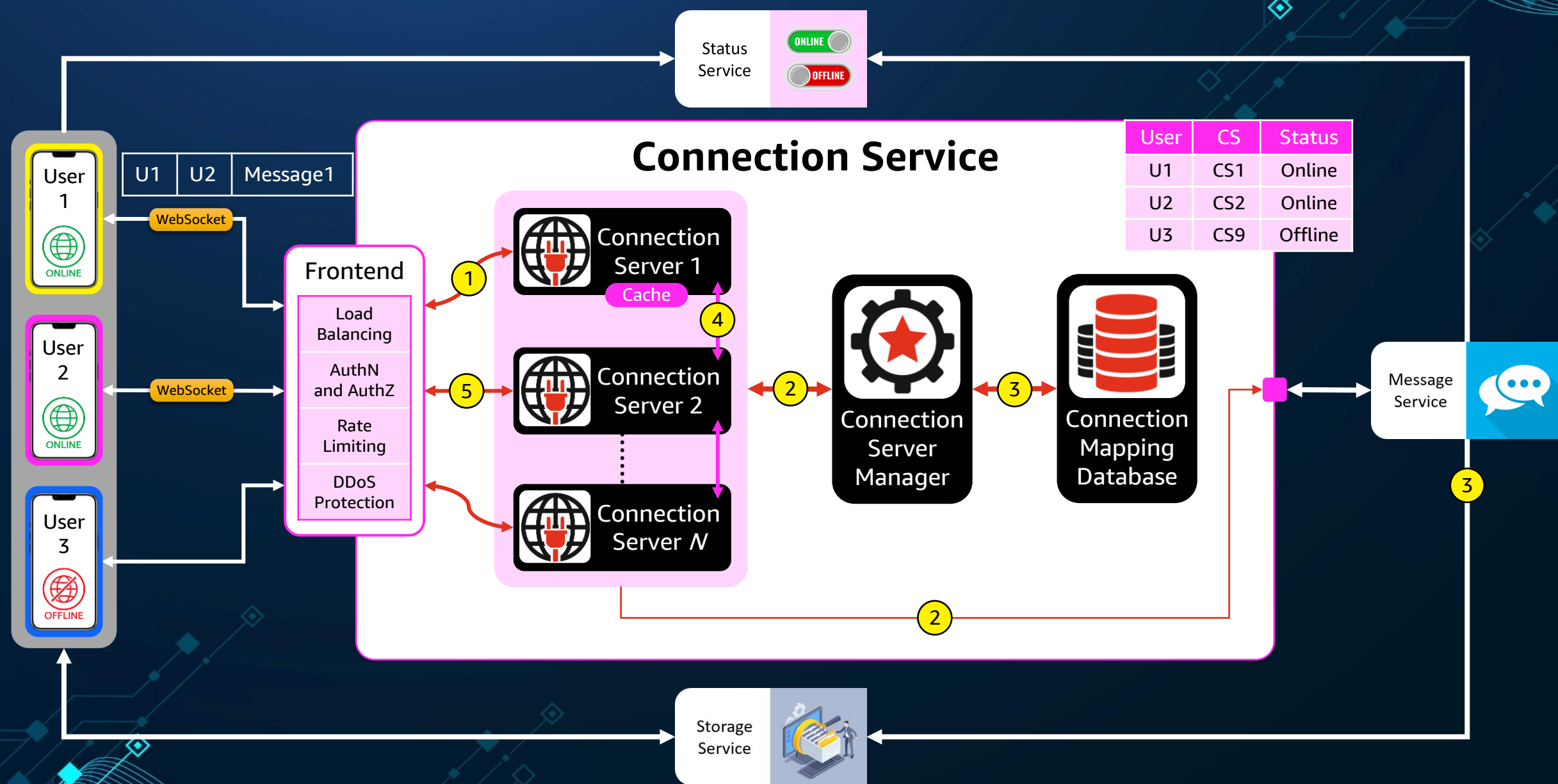- Connection Handler
- Messaging Module
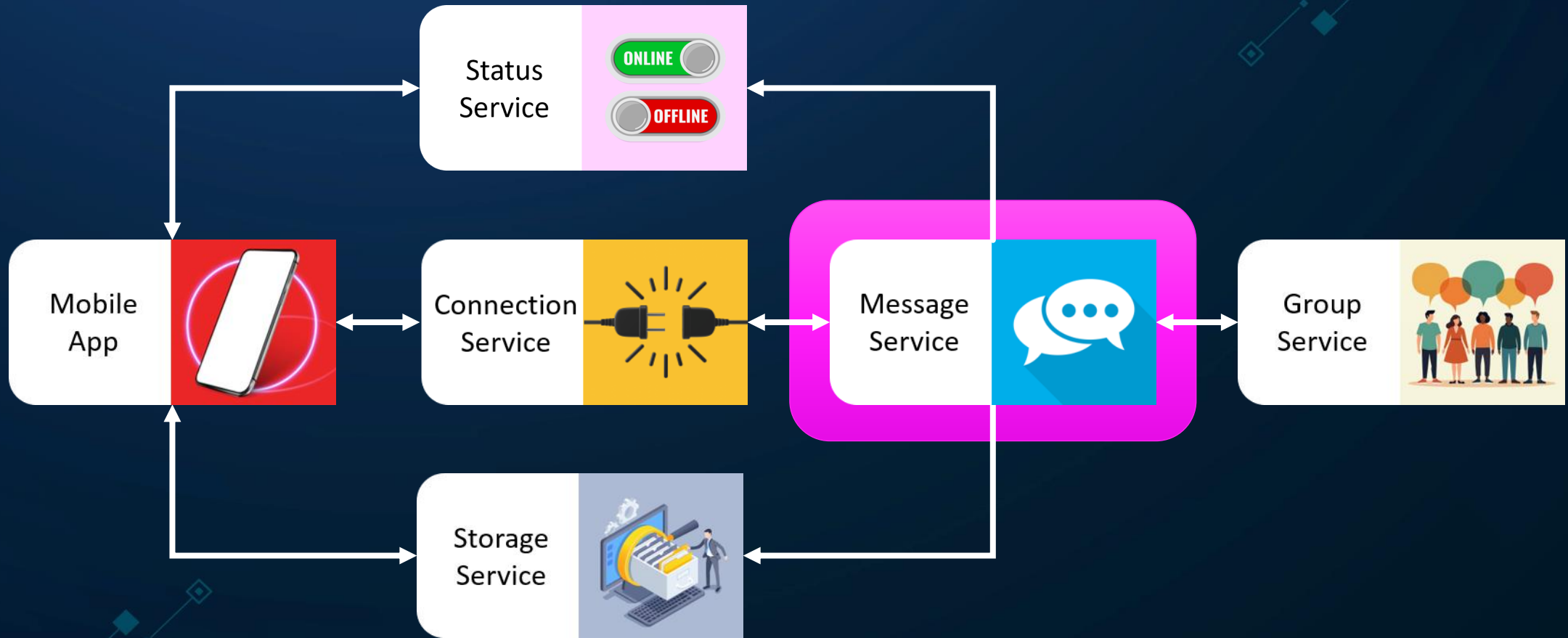- Local Storage & Sync
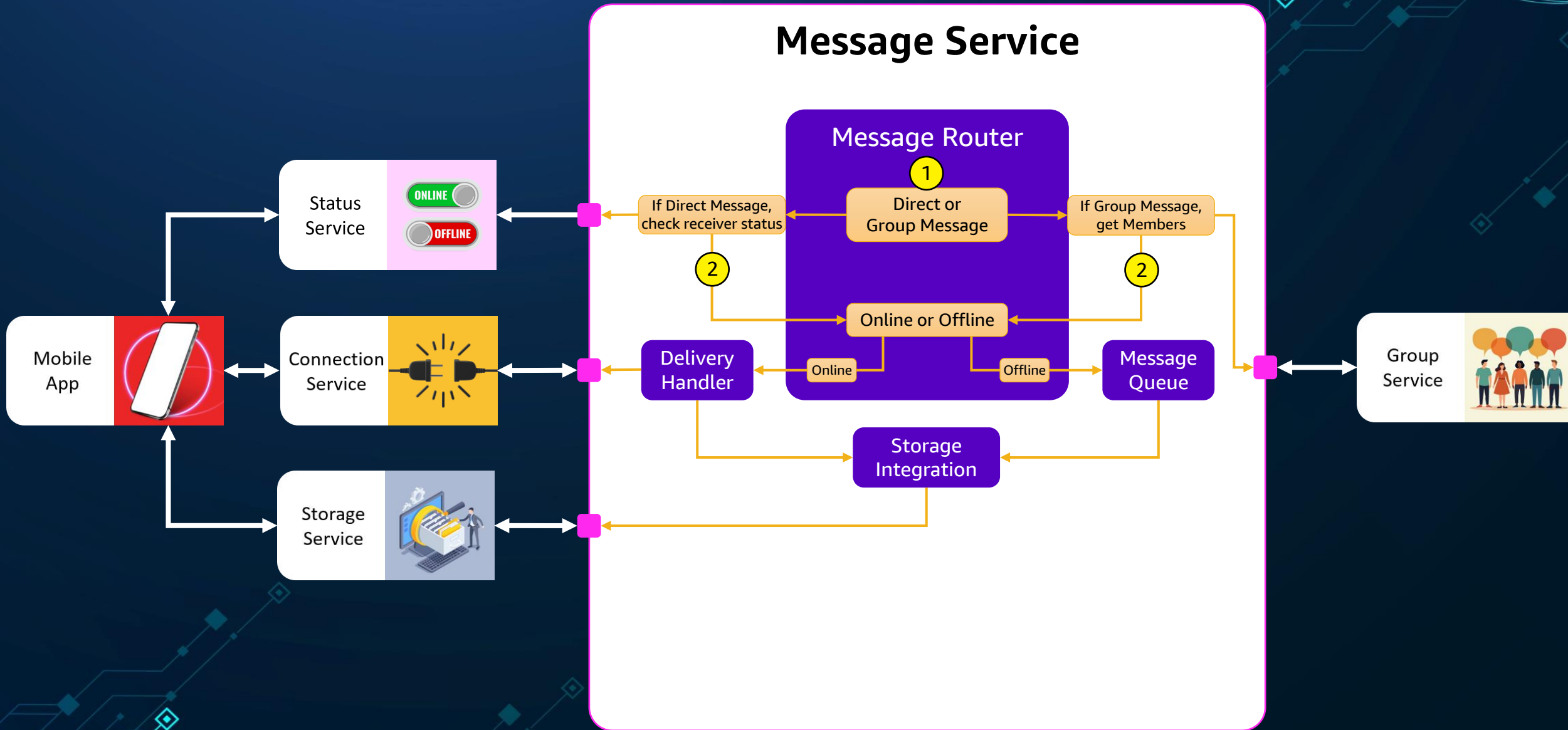- Notification Handler

# Connection Service

# Connection Service

# Message Service

# Message Service



**Message Service**

**Message Router**

Direct or Group Message

If Direct Message, check receiver status

If Group Message, get Members

Online or Offline

Online

Offline

Delivery Handler

Message Queue

Storage Integration

Status Service

ONLINE
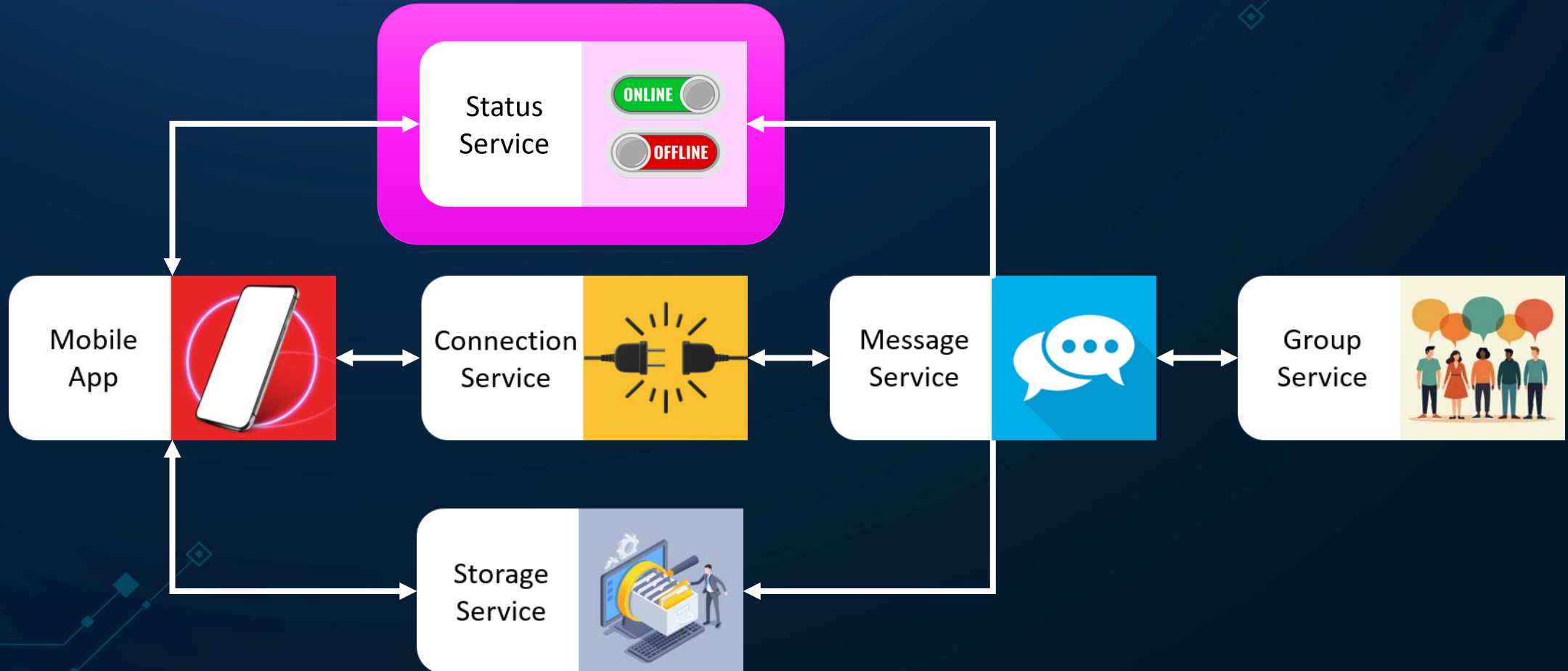
OFFLINE

Connection Service
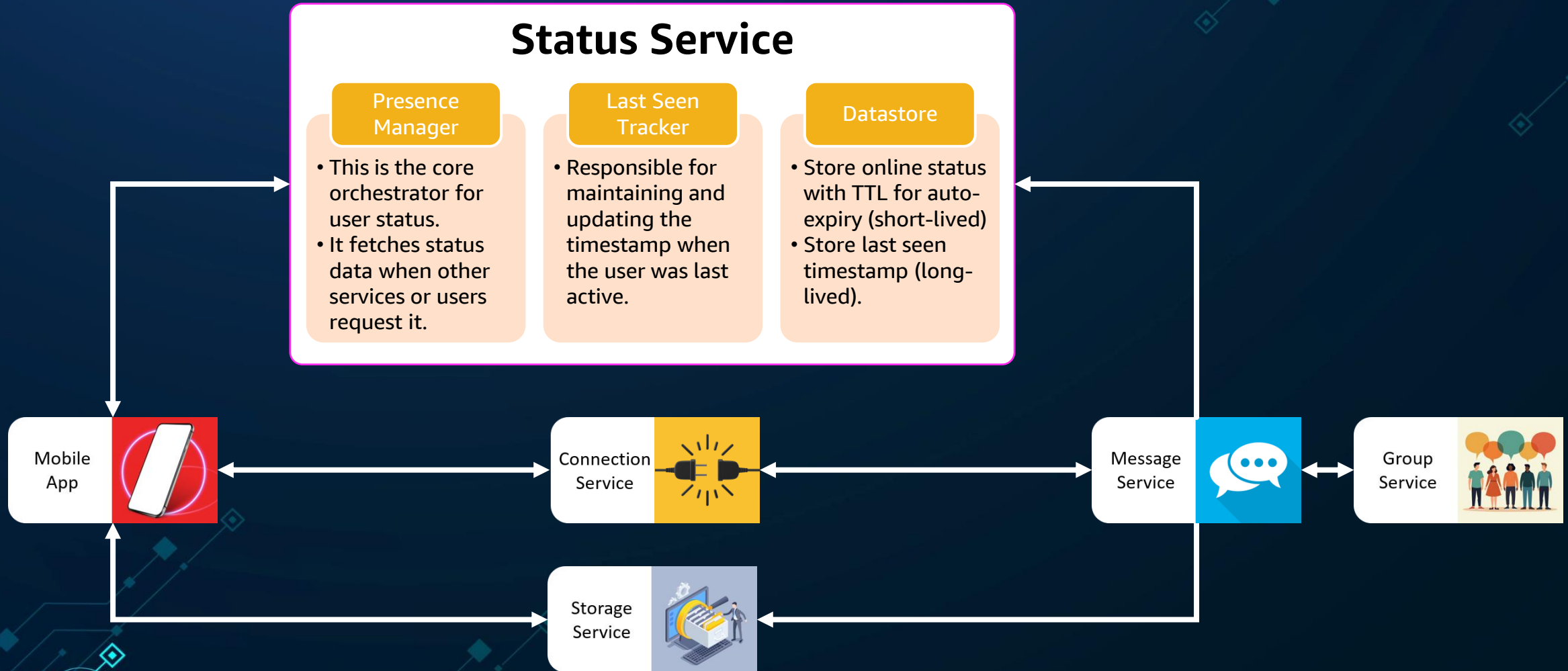
Storage Service

Mobile App

Group Service

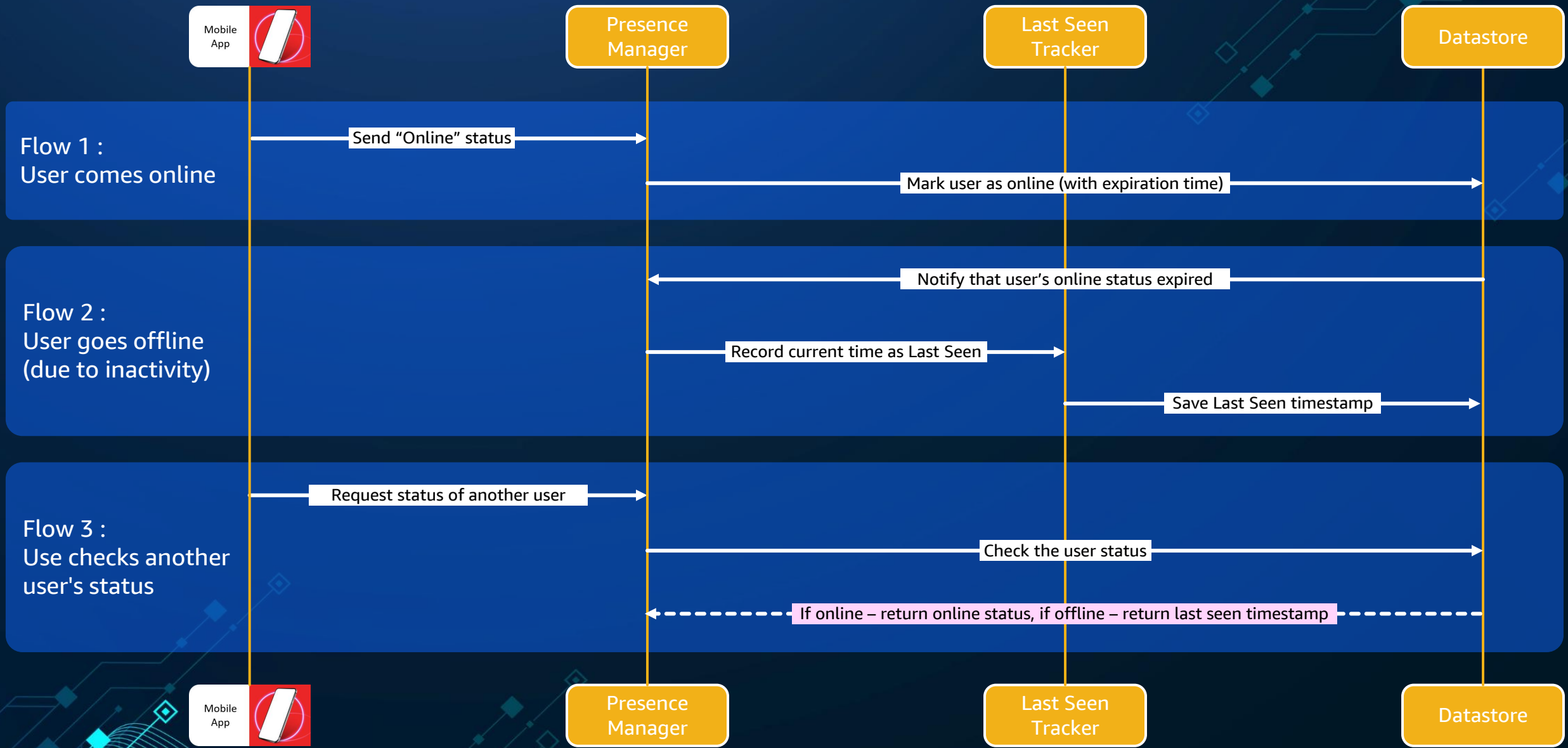# Status Service

- Track user presence (online/offline) and last seen time.

# Status Service

- Track user presence (online/offline) and last seen time.

## Status Service

### Presence Manager

- This is the core orchestrator for user status.
- It fetches status data when other services or users request it.

### Last Seen Tracker

- Responsible for maintaining and updating the timestamp when the user was last active.

### Datastore

- Store online status with TTL for auto-expiry (short-lived)
- Store last seen timestamp (long-lived).

Mobile App

Connection Service

Message Service

Group Service

Storage Service

# Status Service



Flow 1 :
User comes online

- Mobile App → Presence Manager: Send "Online" status
- Presence Manager → Datastore: Mark user as online (with expiration time)

Flow 2 :
User goes offline
(due to inactivity)

- Datastore → Presence Manager: Notify that user's online status expired
- Presence Manager → Last Seen Tracker: Record current time as Last Seen
- Last Seen Tracker → Datastore: Save Last Seen timestamp

Flow 3 :
Use checks another
user's status

- Mobile App → Presence Manager: Request status of another user
- Presence Manager → Datastore: Check the user status
- Datastore ⇠ Mobile App: If online – return online status, if offline – return last seen timestamp
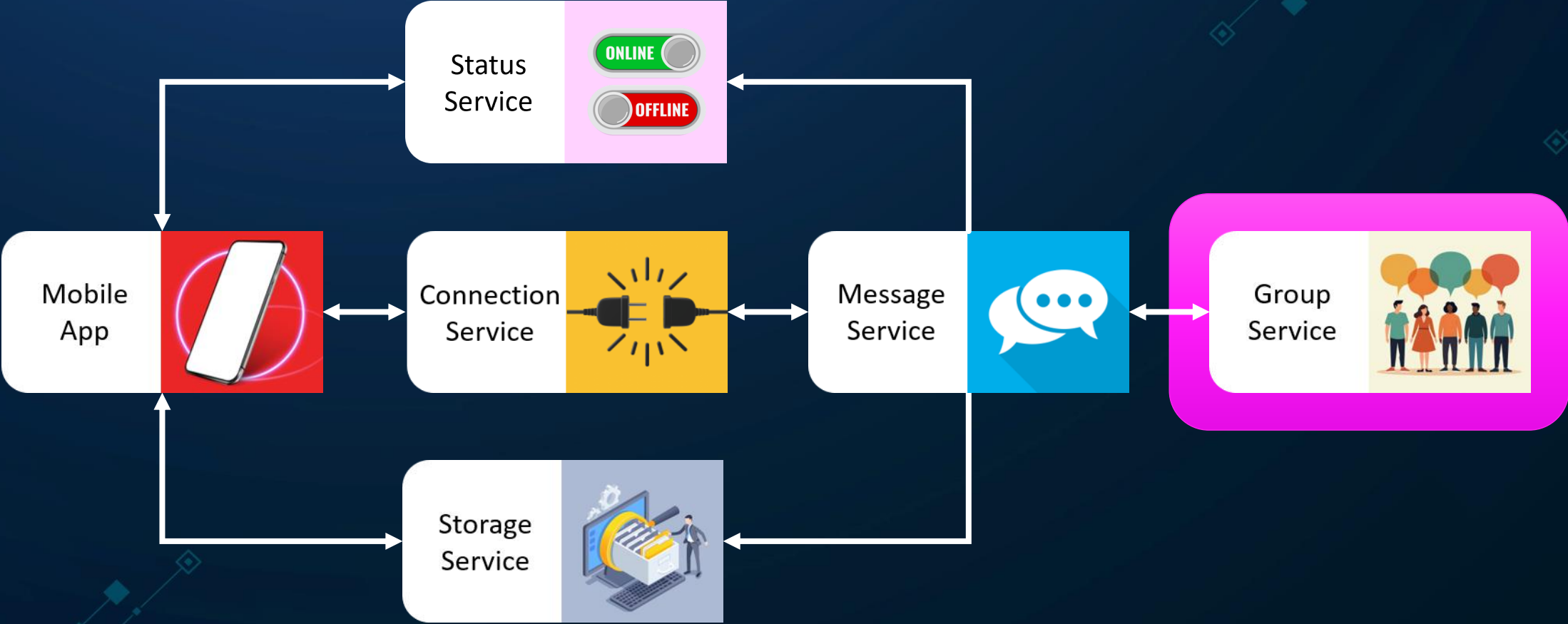
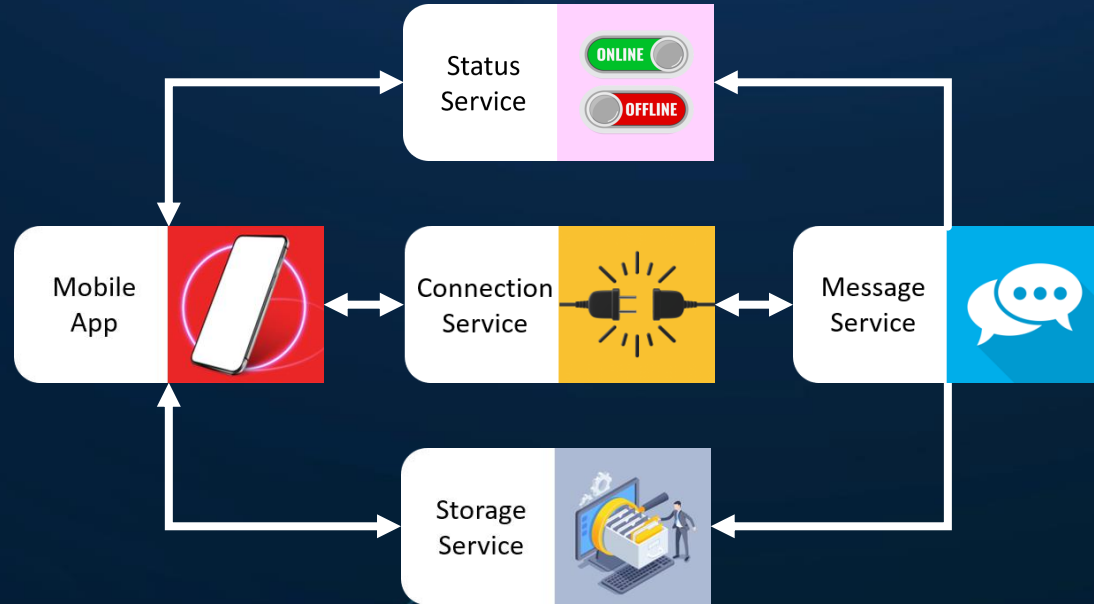Mobile App     Presence Manager     Last Seen Tracker     Datastore

# Group Service
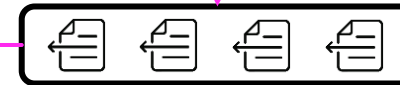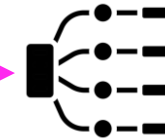
# Group Service



## Group Service

- Fetch group members from DB
- Generate delivery tasks
- Push to central queue

**Fan Out Engine**

- Store group info –
- Maintain member lists
- Support permissions/ admin logic

**Group Membership Database**

**Central Message Queue**

- Buffer delivery jobs
- Maintain order & durability
- Allow retries

Status Service

ONLINE
OFFLINE

Mobile App

Connection Service

Message Service

Storage Service

# Fan out design – Analogy

- Want to share a party invitation to your colleagues

| | | |
|---|---|---|
| Send individual letters to each person | Post the invitation once on a notice board | Cards for close friends, board for everyone else |
| Everyone gets it quickly, but it's resource-intensive | Efficient and scalable, but delivery timing depends on the reader | Mix of speed (for VIPs) and scale (for general public) |
| Fan-out on Write | Fan-out on Read | Hybrid Strategy |

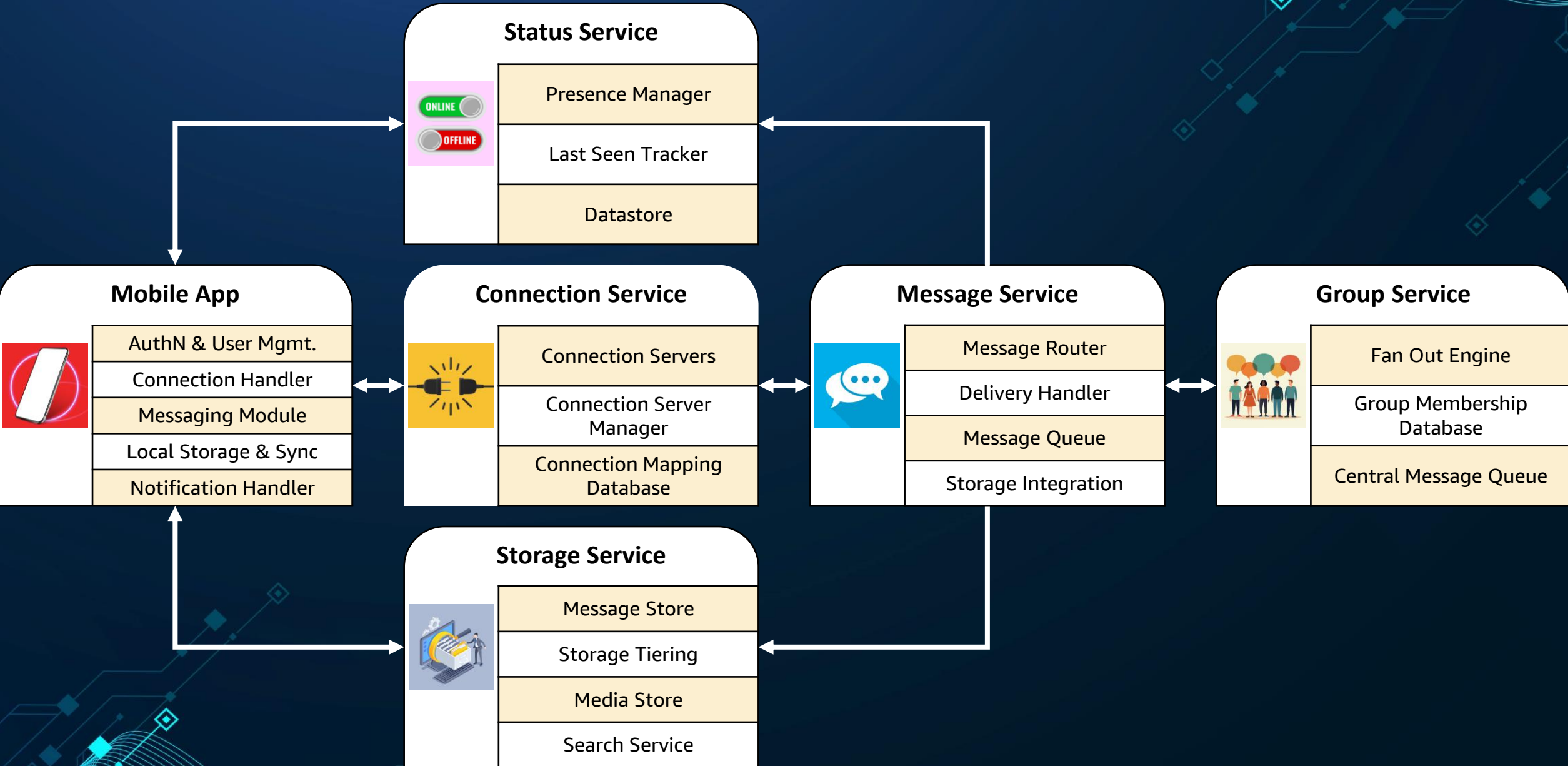| Approach | How It Works? | Pros | Cons | Best Suited For |
|----------|---------------|------|------|-----------------|
| **Fan-out on Write** | Message is written to each group member's queue or inbox at send time | • Low latency delivery<br>• Simple read path<br>• Easy offline support | • High write amplification (1 → N)<br>• Costly for large groups-<br>• Duplicate storage<br>• Complex retries on failure | • Small to medium groups (<1,000 members)<br>• Private group chats |
| **Fan-out on Read** | Message is stored once; delivered/fetched by users when they come online or request it | • Write once, read many<br>• Storage-efficient<br>• Ideal for large groups<br>• Scales to millions of members | • Higher read latency<br>• Must track user read pointers<br>• Harder for real-time delivery<br>• More complex read logic | • Large/public groups (>10K members)<br>• Broadcast channels |
| **Hybrid (Tiered)** | Uses fan-out on write for small groups, fan-out on read for large groups | • Flexible and scalable<br>• Optimizes both read and write<br>• Supports varied group types and use cases | • Architecturally complex<br>• Must manage two delivery models<br>• Higher dev and testing effort | • Platforms like WhatsApp<br>• Tiered user experiences (e.g., premium vs. free) |

# Storage Service

# Storage Service



**Status Service**

**Mobile App**

**Connection Service**

**Message Service**

**Group Service**

## Storage Service

### Message Store
- Stores chat messages along with metadata, including delivery status.

### Storage Tiering
- Archives less frequently accessed data to cost-effective cold storage.

### Media Store
- Handles storage of all media files (images, videos, documents, audio, etc.).

### Search Service
- Enables efficient message lookups and supports filters and queries.

# Storage Tiering

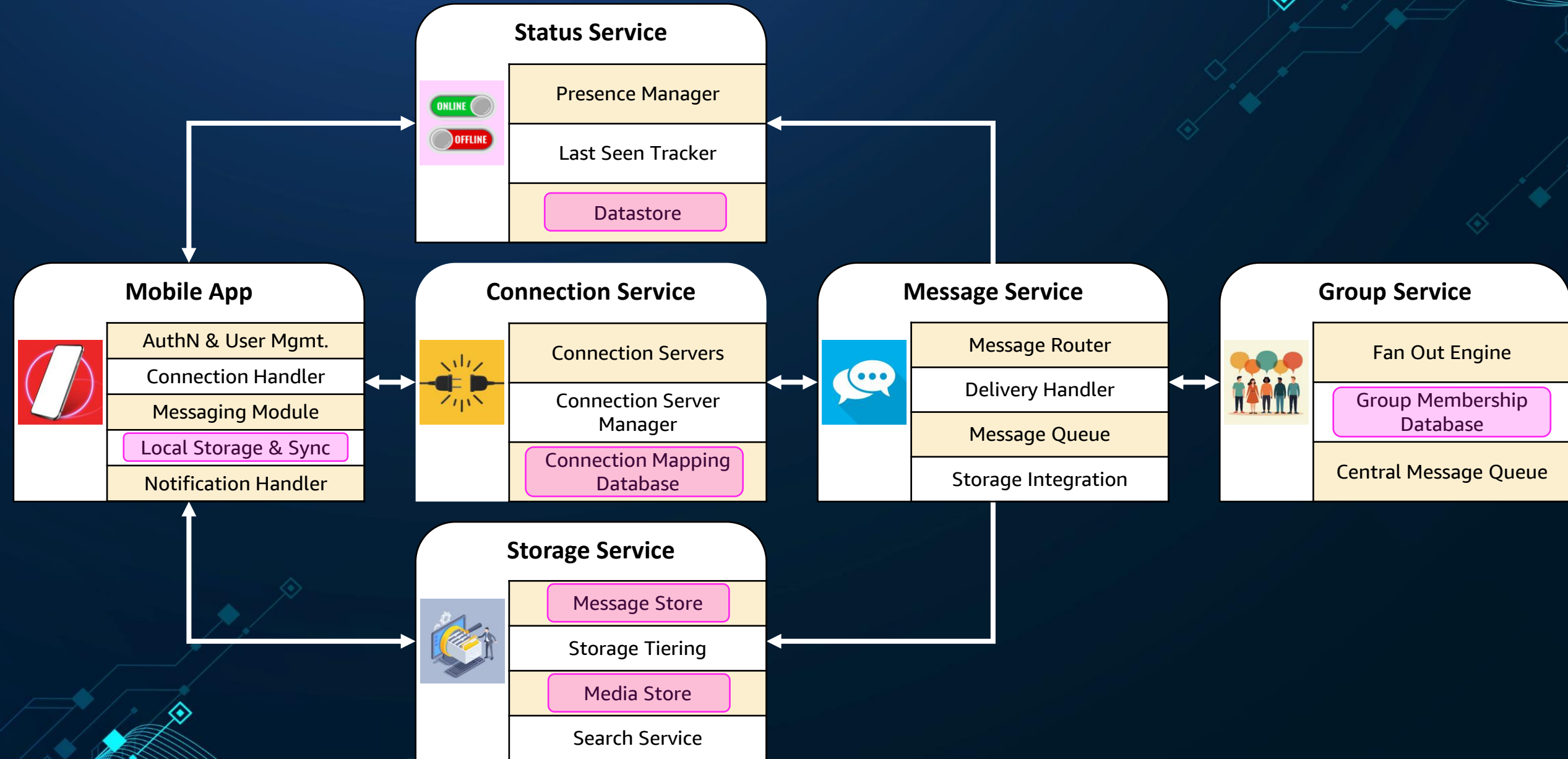| Tier | Purpose | Storage Type (Example) |
|------|---------|------------------------|
| Hot Storage | For active, recently accessed media | High-speed object store (e.g., S3 Standard) |
| Warm Storage | For moderately accessed media | Cheaper but slower object storage (e.g., S3 Infrequent Access) |
| Cold Storage | For archival or rarely accessed data | Archival tier (e.g., S3 Glacier) |

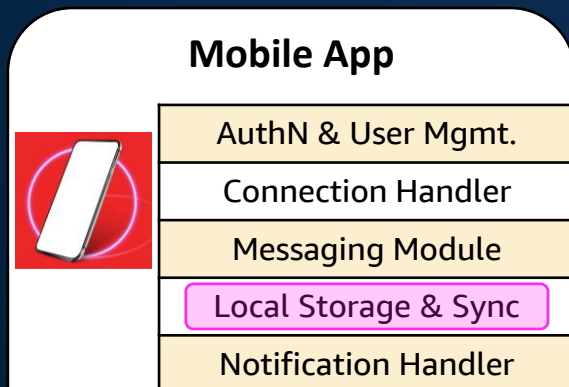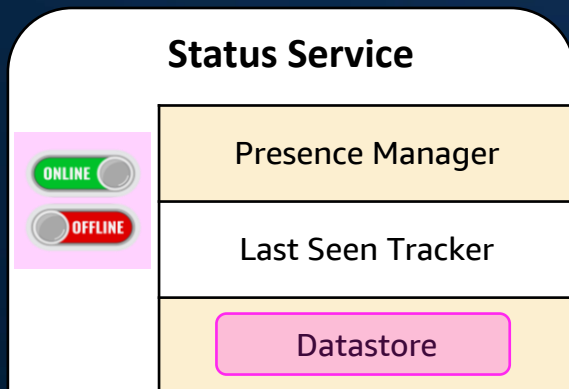| Criteria | Why It Matters? |
|----------|-----------------|
| Last Access Time | Most common — if media hasn't been accessed in X days, move it. |
| Message Age | Move media older than N days/weeks/months regardless of access. |
| Message Read Status | Media in messages that are read by all group members may be cold. |
| User Activity | If sender and receiver are inactive, media can be tiered down. |
| Content Type | Larger video files may move to cold faster than small images. |
| Group Size & Engagement | Media from inactive or low-engagement groups can be cold-tiered. |

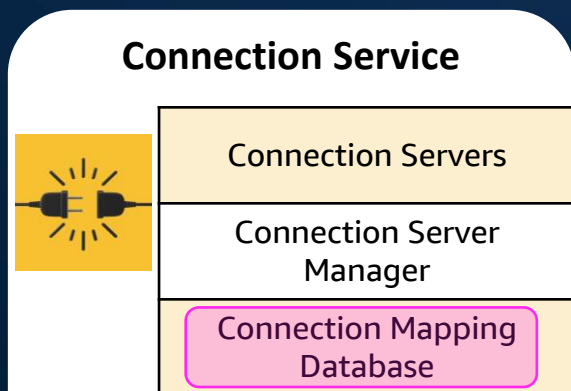# Complete architecture

# Which storage?

| Sub-component | Mobile App > Local Storage |
|---|---|
| Purpose | • Store chat history<br>• Message metadata<br>• User preferences locally |
| Consideration | • Structured, relational data<br>• Needs offline access<br>• Must support search<br>• Pagination<br>• Joins |
| Relational or Non-relational | **Relational**<br>• Lightweight and embeddable on mobile devices<br>• Supports ACID compliance for reliability<br>• Efficient for structured, tabular data (messages, chats, users)<br>• Optimized for complex local queries and indexes<br>• E.g. - SQLite |

**Mobile App**

- AuthN & User Mgmt.
- Connection Handler
- Messaging Module
- Local Storage & Sync
- Notification Handler

## Status Service

| | |
|---|---|
| ONLINE / OFFLINE | **Status Service** |
| | Presence Manager |
| | Last Seen Tracker |
| | Datastore |

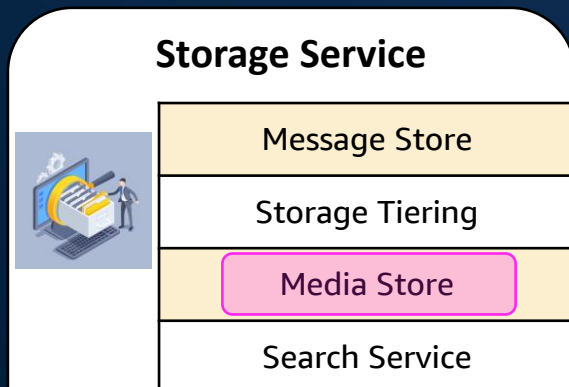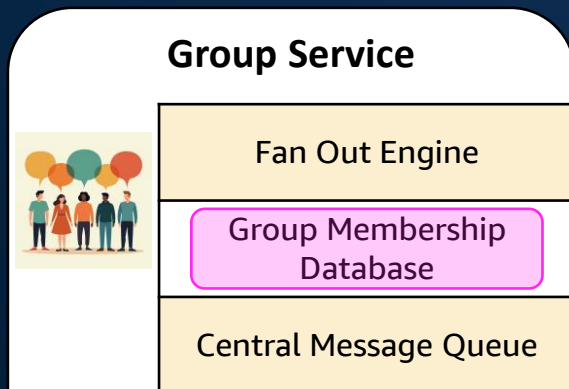| | |
|---|---|
| **Sub-component** | Status Service > Datastore |
| **Purpose** | • Track online/offline status<br>• Typing indicators<br>• Last seen |
| **Consideration** | • Ephemeral and fast-changing<br>• Needs very low latency<br>• Time-based auto-expiry needed – Time-To-Live (TTL) |
| **Relational or Non-relational** | **Non-relational**<br>• Ideal for high-frequency updates and quick lookups<br>• Supports TTL for auto-expiry of status<br>• In-memory option enables real-time speed<br>• Horizontally scalable with predictable performance<br>• No schema rigidity — easy to evolve status format<br>• E.g., Redis, DynamoDB |

## Connection Service

- Connection Servers
- Connection Server Manager
- **Connection Mapping Database**

| Sub-component | Connection Service > Connection Mapping Database |
|---|---|
| Purpose | • Store mapping of user ID to active connection server |
| Consideration | • Real-time connection state<br>• High concurrency<br>• Simple key-value access pattern |
| Relational or Non-relational | **Non-relational**<br>• Key-value store perfectly suits user-to-connection mapping<br>• High-throughput and low-latency for billions of users<br>• Built-in TTL support for auto-cleanup of stale connections<br>• Scales horizontally across regions/data centers<br>• No need for relational constraints or joins<br>• E.g., Redis, DynamoDB |

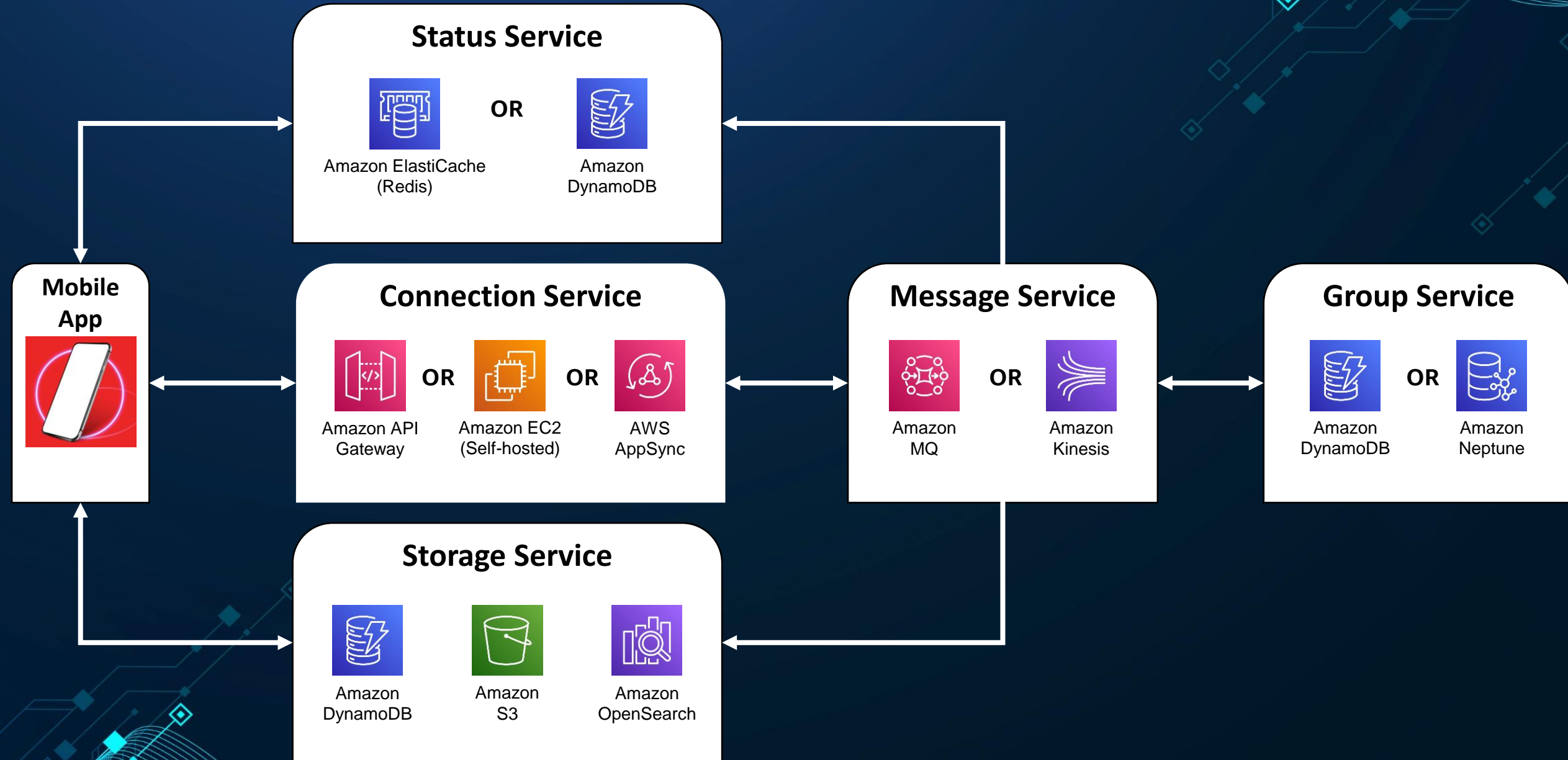| | |
|---|---|
| **Sub-component** | Storage Service > Message Store |
| **Purpose** | • Store and query text messages and metadata (excluding media) |
| **Consideration** | • Massive write throughput<br>• Append-only model<br>• Partitioning by user/chat needed<br>• Flexible schema |
| **Relational or Non-relational** | **Non-relational**<br>• Supports high-velocity writes and sequential reads<br>• Tuned for append-only workloads (no updates)<br>• Horizontal scalability with consistent performance<br>• Schema flexibility as message format evolves<br>• Denormalized design improves read performance for chat history<br>• E.g., Cassandra, DynamoDB |

**Storage Service**

- Message Store
- Archival Manager
- Media Store
- Search Service

| | |
|---|---|
| **Sub-component** | Storage Service > Media Store |
| **Purpose** | • Store media files<br>    • Images, Videos, Document, Audio |
| **Consideration** | • Large binary data<br>• Static after upload<br>• Size limits<br>• Retention limit |
| **Relational or Non-relational** | **Object Storage**<br>• Built to handle unstructured, large binary objects<br>• Content is immutable and served via CDN links<br>• Scalable to petabytes with high availability<br>• Metadata can be indexed in a relational or NoSQL DB<br>• Lifecycle management (e.g., archival, deletion) is built-in<br>• E.g., Amazon S3, GCS, Azure Blob |

**Storage Service**

| Message Store |
|---|
| Storage Tiering |
| Media Store |
| Search Service |

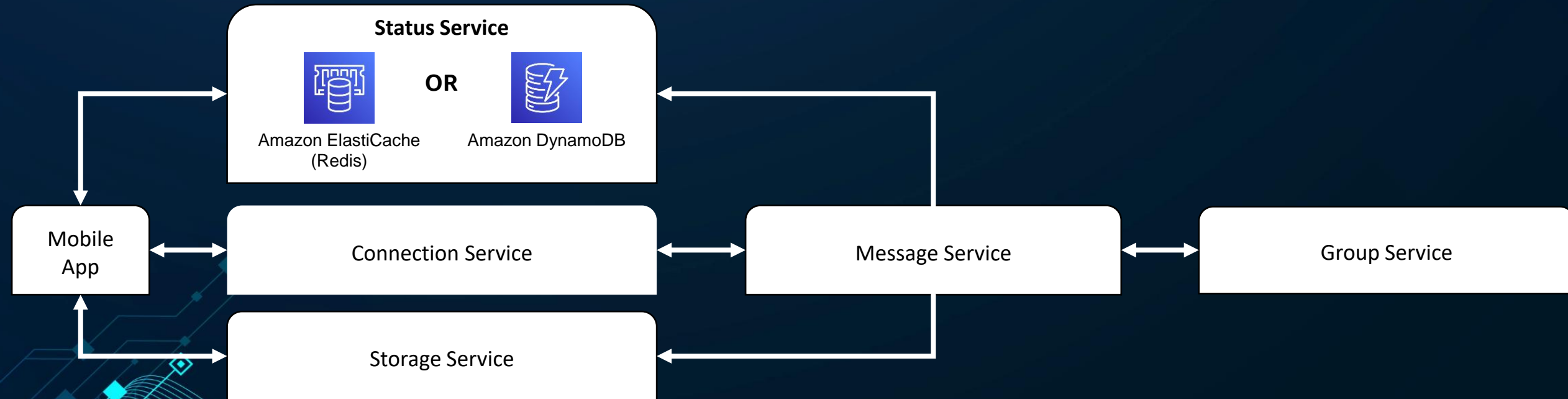| | |
|---|---|
| **Sub-component** | Group Service > Group Membership Database |
| **Purpose** | • Store user-to-group mapping<br>• Roles and permissions |
| **Consideration** | • Many-to-many relationships<br>• High fan-out read performance<br>• Requires fast group membership resolution |
| **Relational or Non-relational** | **Non-relational**<br>• Can handle denormalized lists of group members at scale<br>• Optimized for fast fan-out during group messaging<br>• Eliminates costly joins needed in relational models<br>• Supports flexible role/permission metadata per member<br>• E.g., Cassandra, HBase |

**Group Service**

Fan Out Engine

Group Membership Database

Central Message Queue

# Implementing on AWS

**Status Service**

Amazon ElastiCache (Redis) **OR** Amazon DynamoDB

**Mobile App**

**Connection Service**

Amazon API Gateway **OR** Amazon EC2 (Self-hosted) **OR** AWS AppSync

**Message Service**

Amazon MQ **OR** Amazon Kinesis

**Group Service**

Amazon DynamoDB **OR** Amazon Neptune

**Storage Service**

Amazon DynamoDB   Amazon S3   Amazon OpenSearch

# Pros and Cons – Status Service

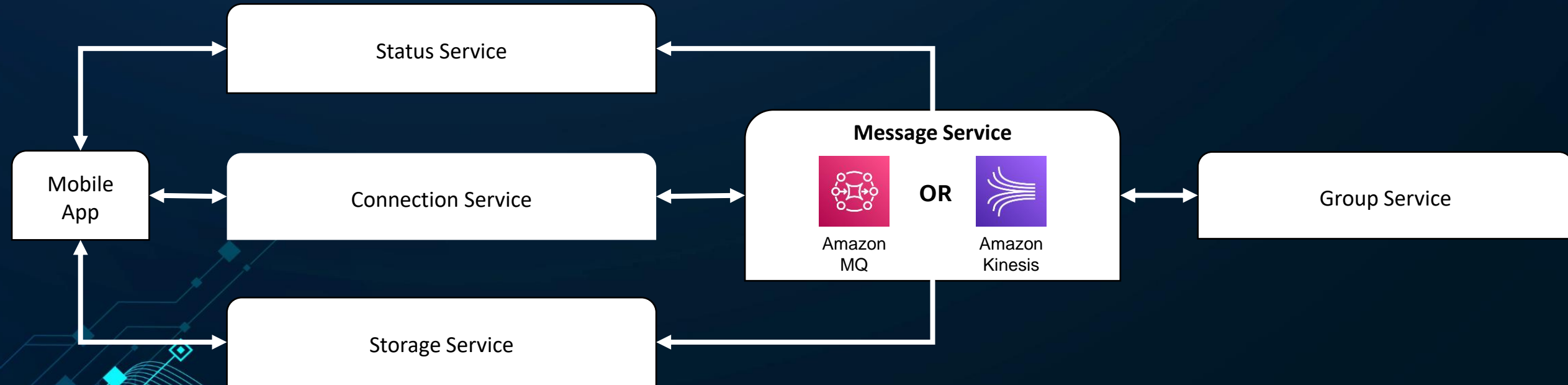| Option | Pros | Cons |
|---|---|---|
| **Amazon ElastiCache (Redis)** | Low latency, Pub/Sub support, built-in TTL for presence. | Not persistent by default; costs increase with scale. |
| **Amazon DynamoDB** | Durable, scalable, supports TTL. | Not real-time; needs DynamoDB Streams + Lambda for triggers. |



**Status Service**

Amazon ElastiCache (Redis)   OR   Amazon DynamoDB

Mobile App

Connection Service

Message Service

Group Service

Storage Service

# Pros and Cons – Connection Service

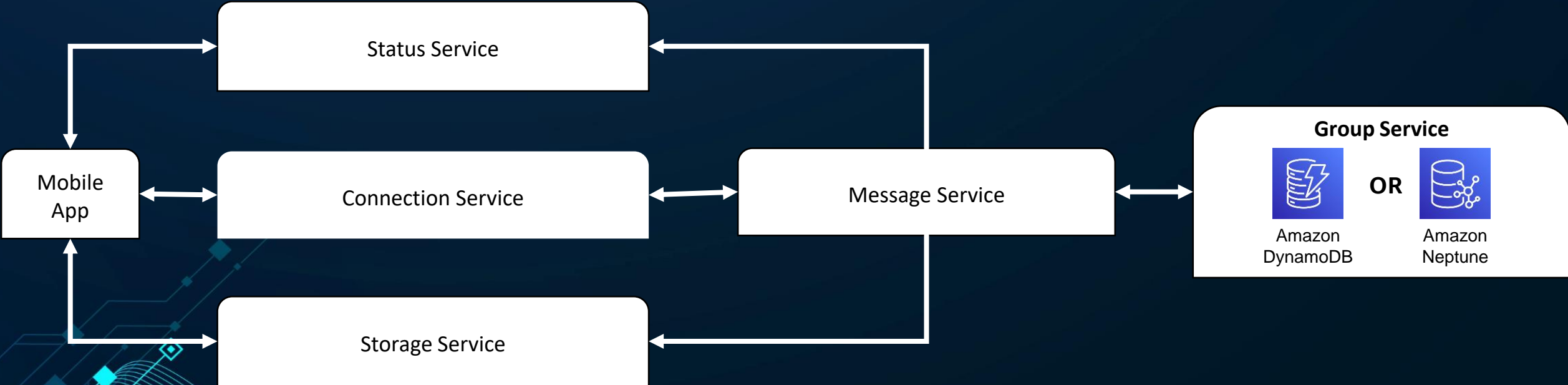| Option | Pros | Cons |
|--------|------|------|
| **Amazon API Gateway** | Native WebSocket support, serverless, scalable. | Limited customization, message size limits. |
| **Amazon EC2 (Self-hosted)** | Full control over connections and protocols. | Requires scaling, monitoring, and maintenance. |
| **AWS AppSync** | Real-time communication, managed. | Might not fit well with highly-custom chat needs. |

# Pros and Cons – Message Service

| Option | Pros | Cons |
| --- | --- | --- |
| **Amazon Kinesis** | High-throughput, real-time stream processing. | Slightly complex to integrate with downstream consumers. |
| **Amazon MQ** | Full-featured message broker. | Heavier, more expensive, managed but not serverless. |

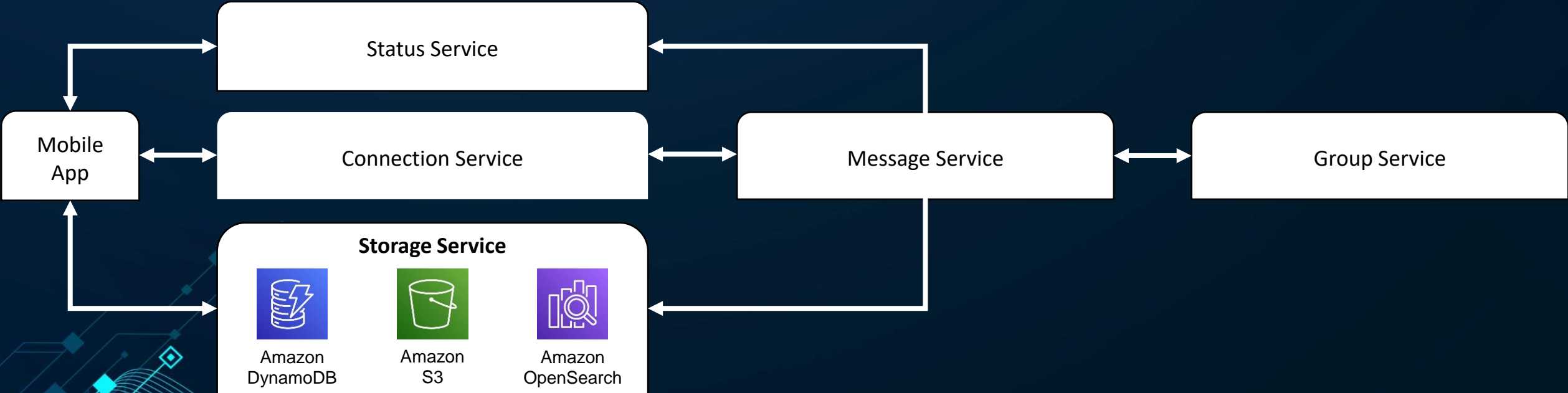# Pros and Cons – Group Service

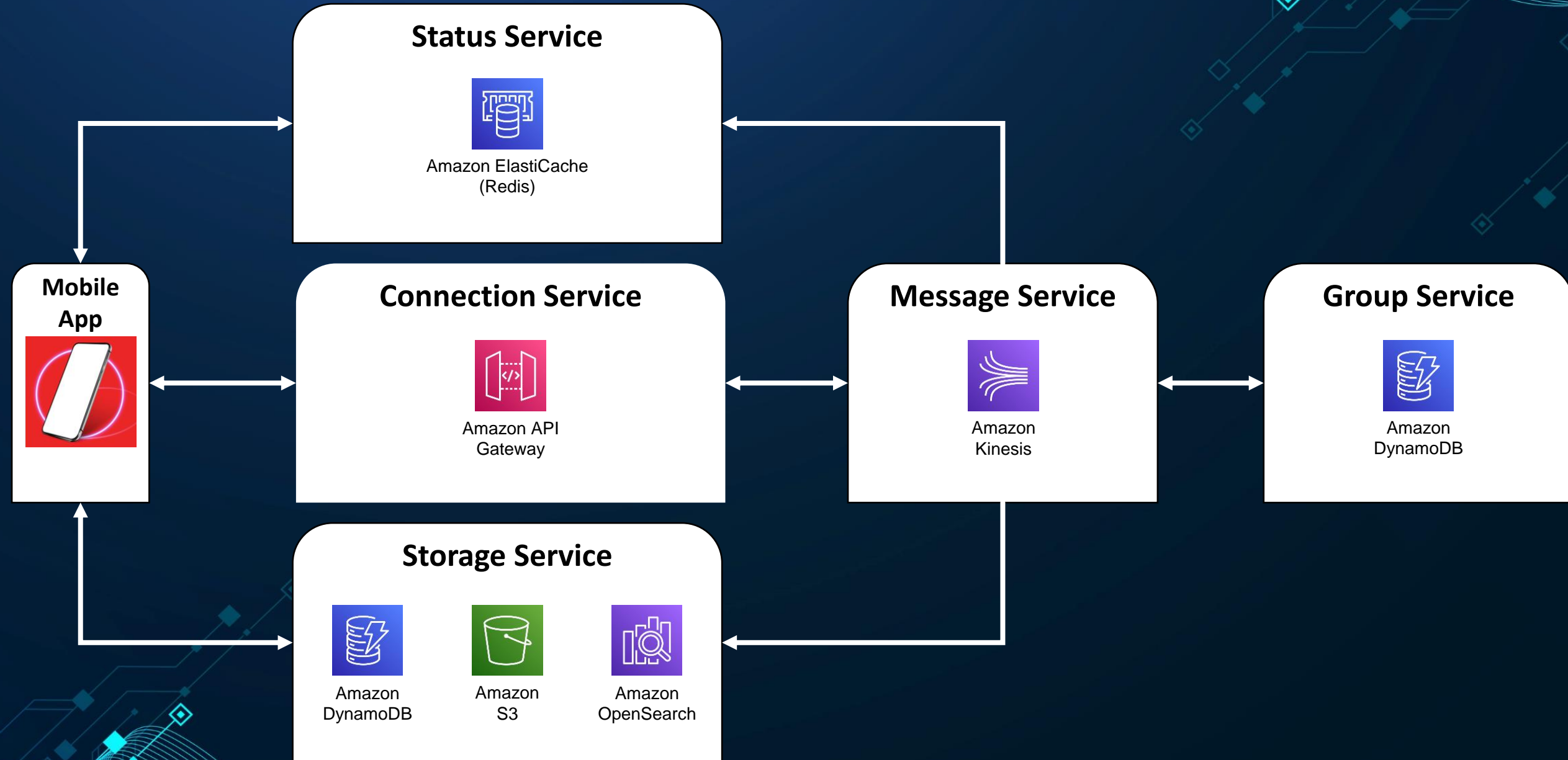| Option | Pros | Cons |
|---|---|---|
| **Amazon DynamoDB** | Scalable, fast key-value access for group/user mapping. | Requires careful key design; no direct relational queries. |
| **Amazon Neptune** | Graph database for complex relationships. | Niche; requires learning curve and integration effort. |

# Pros and Cons – Storage Service

| Subcomponent | Recommended Service(s) |
|---|---|
| **Message Store** | Amazon DynamoDB (with TTL and Streams) |
| **Media Store** | Amazon S3 (Intelligent Tiering, Signed URLs) |
| **Cold Storage Manager** | S3 Lifecycle Policies |
| **Search Capabilities** | Amazon OpenSearch |

# Implementing on AWS

**Status Service**

Amazon ElastiCache (Redis)

**Mobile App**

**Connection Service**

Amazon API Gateway

**Message Service**

Amazon Kinesis

**Group Service**

Amazon DynamoDB

**Storage Service**

Amazon DynamoDB

Amazon S3

Amazon OpenSearch
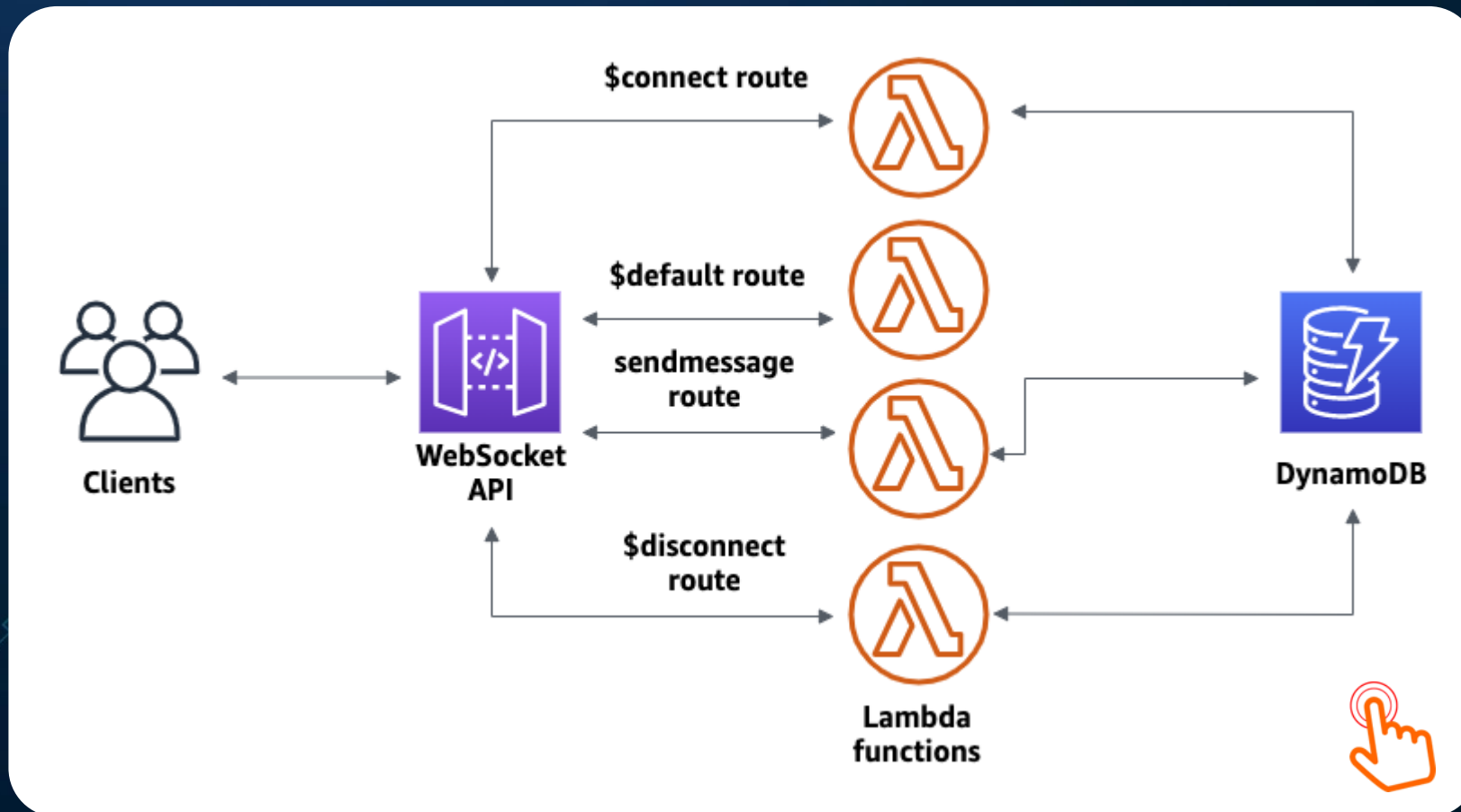
# Tutorial – Create a WebSocket chat app

- Create a WebSocket chat app with a WebSocket API, Lambda and DynamoDB

# Non-functional Requirement



**Performance**

Send and receive messages in real-time

**Scalability**

Billions of users and billions of messages

**Security**

End-to-end encryption of messages

**Reliability**

Store messages until they are delivered

**Availability**

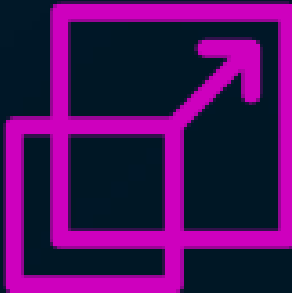Minimal downtime or interruptions

# Performance

- **Goal:** Send and receive messages in real-time

- **Recommendations:**

  - Use persistent connections like WebSockets for bi-directional, low-latency communication.

  - Minimize network hops and reduce the number of intermediary services in the message path.

  - Introduce in-memory caching layers (e.g., Redis-like systems) for user presence, session info, and message routing.

  - Design lightweight message formats (e.g., binary or compacted JSON) to reduce serialization/deserialization overhead.

  - Push vs Pull: Prefer push-based message delivery over polling to reduce latency and resource usage.

  - Load balance traffic across multiple message handlers and connection servers.

# Scalability

- **Goal:** Billions of users and billions of messages

- **Recommendations:**

  - Use stateless microservices wherever possible so they can scale horizontally.

  - Introduce message queues or streams to decouple services (e.g., routing, storage, delivery).

  - Partition data by user, region, or chat to avoid bottlenecks and hotspots.

  - Use elastic storage and compute layers that can automatically expand with traffic.

  - Auto-scale connection and processing layers based on metrics like CPU, memory, queue depth, or user activity.

# Security

- **Goal:** End-to-end encryption of messages

- **Recommendations:**

  - Implement end-to-end encryption (E2EE) so only sender and receiver can read messages—ensure encryption is handled at the client.

  - Use strong user authentication with token-based mechanisms (e.g., OAuth2 or JWT).

  - Encrypt all data in transit using HTTPS/TLS and encrypt at rest using managed keys.

  - Isolate sensitive components using private networks or access-controlled environments.

  - Apply least-privilege access controls for both services and users.

  - Log and monitor all access events for audit and anomaly detection.

# Reliability

- **Goal:** Store messages until they are delivered

- **Recommendations:**

  - Persist messages before delivery to a durable store to avoid data loss during crashes.

  - Implement delivery acknowledgements and retries with exponential backoff.

  - Design for idempotency: duplicate messages or retries should not cause unintended side effects.

  - Use message queues with dead-letter support to handle failed deliveries gracefully.

  - Add health checks and service monitoring for automatic failure detection and failover.

  - Redundancy at all layers: Have multiple replicas of critical services and data stores.
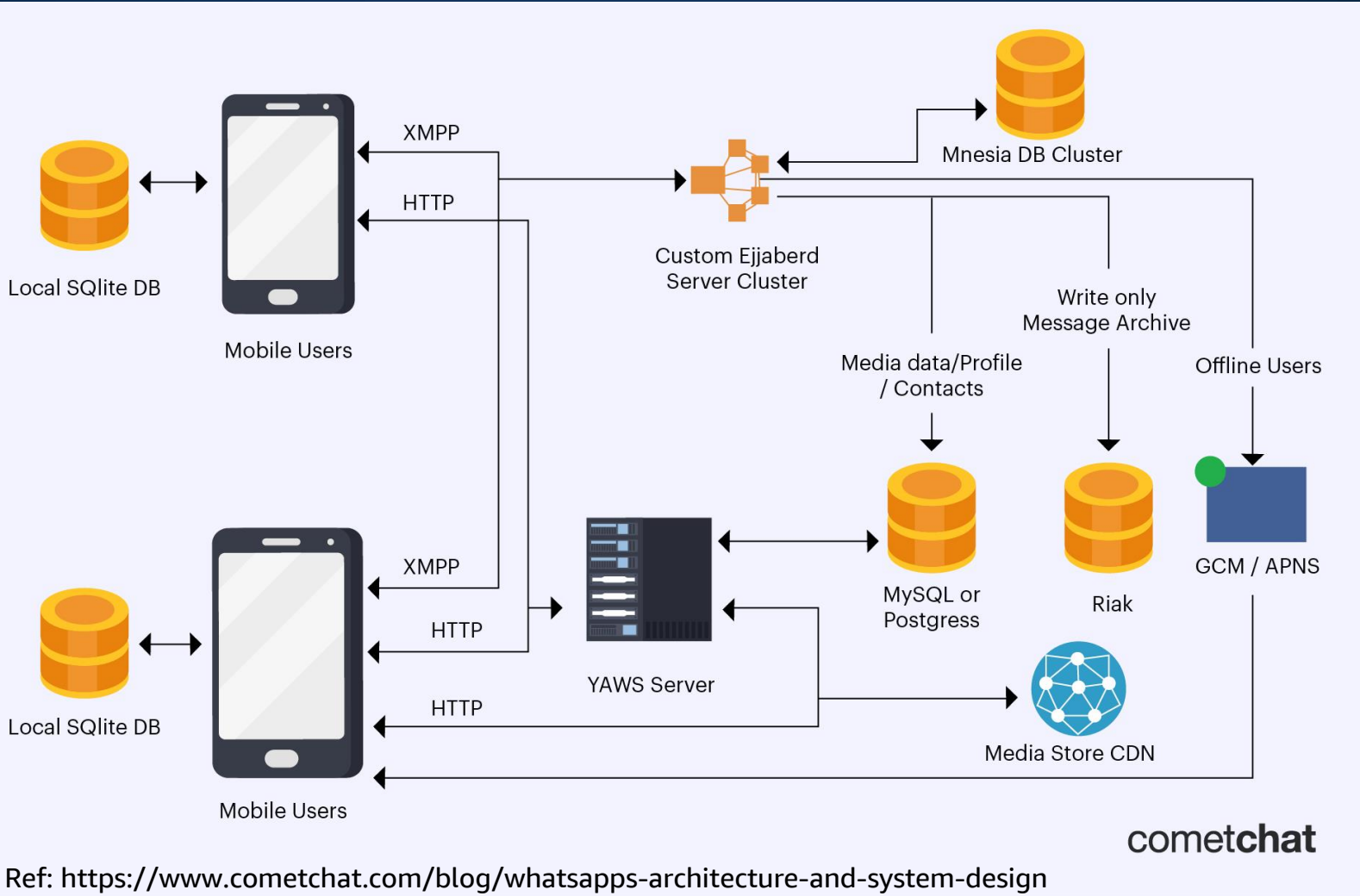
# Availability

- Goal: Minimal downtime or interruptions

- Recommendations:

  - Deploy services across multiple zones or regions to tolerate infrastructure failures.

  - Use load balancers and traffic routers to distribute traffic and handle node failures.

  - Gracefully degrade features: if group messaging fails, allow P2P to continue; if media upload is slow, queue for retry.

  - Implement self-healing infrastructure: automatically restart failed services or reroute traffic.

  - Regularly back up data and validate recovery processes for disaster recovery readiness.

  - Monitor availability with real-time alerts, dashboards, and SLO tracking.

# WhatsApp Architecture – based on publicly shared information

• No officially published architecture



Ref: https://www.cometchat.com/blog/whatsapps-architecture-and-system-design

| Programming Languages | • Erlang |
|---|---|
| Media-related components | • C++ |
| Database | • Mnesia (Erlang's distributed DB)<br>• Later scaled with MySQL (for long-term storage)<br>• RocksDB (for high-speed access).<br>• Riak (for media storage) |
| Messaging Protocol | • FunXMPP |
| Web Server | • YAWS (Yet Another Web Server) |
| End-to-end encryption | • Signal Protocol |
| Infrastructure | • Initially using FreeBSD, Now likely hosted within Meta's global data centers |

| Aspect | WhatsApp | Telegram | WeChat | Facebook Messenger |
|---|---|---|---|---|
| Focus Area | Privacy, simplicity, end-to-end communication | Speed, openness, cloud sync, developer-friendly APIs | Ecosystem of services, mini-apps, payments | Social interaction, integration with FB/Instagram/Threads |
| Unique Architecture Trait | Mobile-first, peer-to-peer encryption core | Cloud-first, centralized but secure transport layer | All-in-one "super app" model (chat, pay, services, games) | Deep integration with Facebook ecosystem |
| Message Protocol | Custom XMPP-based protocol + proprietary extensions | MTProto (Telegram's own protocol, optimized for speed & security) | Custom protocol (based on Tencent's internal standards) | MQTT (lightweight pub/sub messaging protocol) |
| Data Storage | Mostly on-device, only metadata stored on servers | Messages stored in cloud (server-side history) | Server-side storage for chat, mini-programs, payments | Cloud-based storage, integrated with Facebook infra |
| Message Encryption | End-to-End by default (Signal Protocol) | Optional End-to-End (Secret Chats only) | Encrypted during transmission; not E2E by default | E2E only for Secret Conversations (not default) |
| Scalability Approach | Sharded Erlang clusters; client-heavy | Stateless backend, API-centric, CDN for media | Monolithic + Microservices for super app ecosystem | Microservices, large-scale sharding, FB global infra |
| Account Identity | Phone-number based | Phone optional; uses username model | Phone-number based + WeChat ID | Facebook account-based |
| Multi-Device Support | Recently added (client still primary source) | Built-in from start; true cloud sync | Yes; supports multiple devices natively | Fully supported |

# Incorporate following in design

- Authentication
- Encryption
- User registration
- User profiles
- Mobile app
- Audio/Video call
- Block lists
- Multiple device support