

# System Design Calculations Cheat Sheet

Formulas, typical values, and worked examples for fast interview math

## Five Essential Calculations

### 1) REQUESTS PER SECOND (RPS)

**Formula:**  $(DAU \times actions\_per\_user \times peak\_ratio) \div 86,400$

- Typical **peak\_ratio**: 2-3x average (can be 5-10x for events).
- Always sanity-check against product behavior (bursty vs steady).

**Example:** 100M DAU, 10 actions/day, 3x peak  $\rightarrow (100,000,000 \times 10 \times 3) \div 86,400 \approx 34,722$  RPS (~35K).

### 2) STORAGE REQUIREMENTS

**Formula:**  $entities \times size\_per\_entity \times growth\_years \times replication\_factor$

- Typical **replication\_factor**: 2-3x (storage systems often keep multiple copies).
- Separate **hot** vs **cold** data if possible (cost and access patterns differ).

**Example:** 1B photos, 2MB each, 3 years, 3x replication  $\rightarrow 1,000,000,000 \times 2MB \times 3 \times 3 \approx 18,000,000,000$  MB  $\approx 18$  PB.

### 3) BANDWIDTH

**Formula:**  $RPS \times avg\_request\_size$  (compute read and write separately)

- Typical request size: 1KB-1MB depending on content type.
- Remember: 1 byte = 8 bits. Convert to Gbps when talking networks.

**Example:** 10K RPS  $\times$  100KB = 1,000,000 KB/s  $\approx$  1 GB/s  $\approx$  8 Gbps.

### 4) CACHE SIZING

**Formula:**  $total\_data \times 0.20$  (80/20 hot set)  $\times$  1.5 (overhead)

- Typical cache hit rate target: 70-90% for well-designed systems.
- Overhead accounts for metadata, eviction policies, and fragmentation.

**Example:** 100GB total data  $\rightarrow 100 \times 0.20 \times 1.5 = 30$ GB cache needed.

### 5) DATABASE CAPACITY

**Rule of thumb:** If calculated RPS exceeds a single instance capacity, plan for replication, caching, and eventually sharding.

Store Type	Typical Single Instance Throughput (very workload-dependent)
SQL (e.g., MySQL/PostgreSQL)	1K-10K QPS
NoSQL (e.g., Cassandra/Dynamo-style)	10K-100K+ QPS

## Reality Check Benchmarks

Component	Benchmark (rough)	Why it matters
Single app server	1,000-5,000 RPS (simple ops)	If you need 100K RPS, you will need many app servers + load balancer
Redis	100,000+ ops/sec	Great for hot data and rate-limiting; reduces DB load.
PostgreSQL	5,000-15,000 QPS (indexed)	Indexes and query shape dominate; treat as workload-dependent.
CDN	70-90% origin load reduction	For static/media-heavy apps, CDN can cut bandwidth and latency
Most mobile apps	< 100 RPS	Many interview designs are over-engineered; scale should drive cost

## Worked Example: Design Twitter Feed

**Requirements:** 300M DAU, each user loads feed 10 times/day, assume 3x peak traffic.

Step	Calculation	Result (approx)
1) RPS	$300M \times 10 \times 3 \div 86,400$	$\approx 104,167$ RPS (~104K)
2) Storage	$300M \text{ users} \times 100 \text{ tweets avg} \times 280 \text{ chars} \times 3 \text{ years}$	$\approx 25$ TB (order of magnitude)
3) Bandwidth	$104K \text{ RPS} \times 10KB \text{ avg response}$	$\approx 1.0$ GB/s (~8 Gbps)
4) Cache	$25TB \times 0.20 \times 1.5$	$\approx 7.5$ TB cache
5) DB capacity	104K RPS vs single-instance capacity	Exceeds → sharding + replication + caching

**Interpretation:** ~104K RPS is beyond a single database instance. You need to reduce read load (cache + precompute), replicate for reads, and shard when a single primary cannot keep up. Use calculations to justify architecture decisions.