# The Client–Server Architectural Style

A Comprehensive Reference for Request/Reply Distributed Systems

## Contents

**17 Conclusion**                                                                                      **23**

# 1 Overview

The client–server style is a component-and-connector architectural style that structures a system as a set of clients that request services and servers that provide them. This fundamental pattern underlies much of modern distributed computing, from web applications to database systems to enterprise service architectures.

In the client–server style, the division of responsibility is asymmetric: clients initiate interactions by requesting services, while servers wait for requests and respond to them. This asymmetry creates a clear separation of concerns—clients focus on presenting information and capturing user intent, while servers focus on processing requests, managing data, and enforcing business rules.

The style is characterized by its request/reply interaction model. A client sends a request to a server, the server processes the request and produces a result, and the server sends a reply back to the client. This synchronous interaction pattern, while simple conceptually, admits many variations in terms of communication protocols, data formats, error handling, and performance optimization.

## 1.1 Scope and Applicability

The client–server style applies broadly to systems requiring separation between service consumers and service providers. This includes web applications where browsers act as clients to web servers, database applications where application servers act as clients to database servers, enterprise systems where multiple clients share centralized business logic and data, mobile applications communicating with backend services, microservices architectures where services act as both clients and servers, and API-based integrations between systems.

The style is particularly valuable when multiple clients need access to shared resources or services, when centralized control over data or business logic is required, when clients have limited computational resources, when services must be independently deployable and scalable, and when clear boundaries between presentation and processing are beneficial.

## 1.2 Historical Context

The client–server style emerged in the 1980s as an evolution from centralized mainframe computing. Time-sharing systems gave way to networked personal computers accessing shared servers. This transition distributed user interface responsibility to clients while keeping data and business logic centralized.

The style has evolved through several generations. Two-tier architectures connected fat clients directly to database servers. Three-tier architectures introduced application servers between clients and databases. N-tier architectures added further layers for specific concerns. Web-based architectures made browsers the universal thin client. Service-oriented and microservices architectures decomposed servers into independently deployable services.

Understanding this evolution helps architects recognize the trade-offs inherent in different client–server configurations and select appropriate patterns for their context.

## 1.3 Relationship to Other Styles

The client–server style relates to several other architectural styles. It specializes the more general call-return style by adding the asymmetric client and server roles. It can be composed with the

layered style, where each layer acts as a server to the layer above and a client to the layer below. It contrasts with peer-to-peer styles where components have symmetric roles. It provides the interaction pattern used within service-oriented and microservices architectures. It can be implemented using the pipe-and-filter style for request processing pipelines.

Many systems combine client–server interactions with other styles. A microservices system uses client–server between services while employing event-driven patterns for asynchronous communication. A web application uses client–server between browser and server while using the MVC pattern within the server.

# 2    Elements

The client–server style comprises three categories of elements: client components that request services, server components that provide services, and request/reply connectors that enable communication between them.

## 2.1    Client Components

A client is a component that invokes services of a server component. Clients initiate interactions and depend on servers for functionality they do not provide themselves.

### 2.1.1    Types of Client Components

Client components appear in various forms depending on the system context. Thin clients provide primarily user interface functionality, delegating most processing to servers; web browsers exemplify this type. Thick or fat clients contain significant application logic, using servers primarily for data access and shared services. Rich clients balance local processing with server interaction, often providing offline capability. Headless clients are non-interactive components that consume server services, such as batch processors or automated agents. Mobile clients are applications on mobile devices with intermittent connectivity and resource constraints. Other servers acting as clients consume services from other servers, common in multi-tier and microservices architectures.

### 2.1.2    Essential Properties of Client Components

When documenting client components, architects should capture several property categories.

Connectivity properties describe how clients connect to servers. Connection management addresses whether connections are persistent or per-request, connection pooling strategies, and reconnection behavior after failures. Service discovery describes how clients locate servers, whether through static configuration, DNS, service registries, or load balancers.

Resource properties characterize client resource usage. Local processing describes computational work performed on the client. Local storage addresses data cached or stored locally. Network usage describes bandwidth consumption and latency sensitivity.

Behavioral properties describe client interaction patterns. Concurrency addresses whether clients make sequential or concurrent requests. Request patterns characterize typical request sequences and dependencies. Timeout handling describes how clients handle slow or unresponsive servers. Retry policies define how clients retry failed requests.

User interaction properties apply to interactive clients. User interface technology describes the presentation layer implementation. Responsiveness requirements specify acceptable latency for user operations. Offline behavior describes client capability when disconnected from servers.

## 2.2   Server Components

A server is a component that provides services to client components. Servers wait for client requests, process them, and return results.

### 2.2.1   Types of Server Components

Server components vary by the services they provide. Web servers handle HTTP requests, serving static content and routing to application logic. Application servers execute business logic, often hosting application frameworks. Database servers manage persistent data storage and retrieval. API servers expose programmatic interfaces for client consumption. File servers provide network access to file storage. Authentication servers handle identity verification and authorization. Message servers facilitate asynchronous message exchange. Caching servers provide high-speed access to frequently requested data.

### 2.2.2   Essential Properties of Server Components

Port properties describe how servers accept connections. Port types enumerate the interfaces exposed by the server. Concurrency limits specify how many simultaneous client connections or requests are supported. Protocol support identifies communication protocols accepted.

Capacity properties characterize server resources and limits. Throughput describes maximum request processing rate. Latency characterizes response time distribution. Resource consumption describes CPU, memory, and I/O usage patterns. Scalability describes how capacity changes with added resources.

Availability properties address server reliability. Uptime requirements specify required availability levels. Failure modes describe how the server fails and what clients observe. Recovery behavior describes restart and state recovery procedures.

State properties describe server state management. Statefulness indicates whether the server maintains client session state. State persistence addresses how state survives server restarts. State sharing describes how state is shared across server instances.

Security properties address access control and data protection. Authentication requirements specify how clients prove identity. Authorization model describes access control enforcement. Data protection addresses encryption and data handling policies.

## 2.3   Request/Reply Connectors

Request/reply connectors enable communication between clients and servers. A connector has two roles: a request role used by clients to send requests and a reply role used by servers to send responses.

### 2.3.1   Types of Request/Reply Connectors

Connectors vary by communication mechanism and protocol. Local procedure calls connect components within the same process space. Remote procedure calls (RPC) enable calls across process

or machine boundaries with location transparency. HTTP connectors use the web protocol for request/reply communication. Message-based connectors use message queues with request/reply correlation. Database connectors use database protocols like JDBC or ODBC. Custom protocol connectors implement application-specific communication.

### 2.3.2   Essential Properties of Request/Reply Connectors

Communication properties describe the fundamental interaction. Locality indicates whether communication is local (same process), remote (same machine), or networked (different machines). Synchrony distinguishes synchronous (client waits for reply) from asynchronous (client continues, reply delivered later) interaction. Reliability describes delivery guarantees, from best-effort to exactly-once.

Protocol properties describe the communication format. Data format specifies serialization such as JSON, XML, Protocol Buffers, or binary. Message structure defines request and response schemas. Versioning describes how protocol versions are managed.

Security properties address communication protection. Encryption specifies whether data is encrypted in transit using protocols like TLS. Authentication describes how connector endpoints verify identity. Integrity protection addresses detection of message tampering.

Performance properties characterize connector efficiency. Latency describes communication delay. Bandwidth describes data transfer capacity. Connection overhead describes cost of establishing connections.

Error handling properties describe failure behavior. Timeout behavior specifies how long to wait for responses. Error propagation describes how errors are communicated to clients. Retry support describes built-in retry mechanisms.

## 3   Relations

Relations in the client–server style define how elements connect and interact.

### 3.1   Attachment Relation

The *attachment* relation associates component ports with connector roles. Client service-request ports attach to the request role of connectors. Server service-reply ports attach to the reply role of connectors.

### 3.1.1   Properties of Attachment

Cardinality describes how many connections are permitted. One-to-one attachment connects one client to one server through a dedicated connector. Many-to-one attachment connects multiple clients to a single server. One-to-many attachment connects one client to multiple servers, as in load-balanced or replicated scenarios.

Binding time indicates when attachments are established. Static binding is fixed at design or deployment time. Dynamic binding is established at runtime, enabling service discovery and load balancing.

Multiplexing describes how multiple interactions share a connection. Dedicated connections carry one request at a time. Multiplexed connections carry multiple concurrent requests.

## 3.2   Invokes Relation

The *invokes* relation indicates that a client uses a service provided by a server. This relation captures the dependency from client to server.

### 3.2.1   Properties of Invokes

Service identification specifies how the invoked service is named or addressed. Invocation style describes whether invocation is request/reply, one-way, or callback-based. Coupling characterizes the degree of dependency between client and server.

## 3.3   Tier Membership Relation

In tiered client–server systems, the *tier membership* relation assigns components to architectural tiers.

### 3.3.1   Properties of Tier Membership

Tier level identifies the tier such as presentation, business logic, or data. Tier responsibilities describe what functions the tier provides. Inter-tier communication describes allowed communication patterns between tiers.

# 4   Computational Model

The computational model describes how client–server systems execute.

## 4.1   Basic Request/Reply Interaction

The fundamental interaction follows a predictable sequence. The client constructs a request message containing the operation to invoke and necessary parameters. The request is transmitted through the connector to the server. The server receives and parses the request. The server executes the requested operation, potentially accessing other services or data. The server constructs a reply message containing results or error information. The reply is transmitted back through the connector to the client. The client receives and processes the reply.

This basic model assumes synchronous, blocking interaction—the client waits for the reply before proceeding. Many variations exist for different requirements.

## 4.2   Synchronous vs. Asynchronous Interaction

In synchronous interaction, the client blocks waiting for the server's reply. This model is simple to reason about and implement but ties client progress to server responsiveness. Client resources are occupied during the wait. Long-running server operations delay the client. Server unavailability directly blocks the client.

In asynchronous interaction, the client continues execution after sending a request, with the reply delivered later through callbacks, futures, or polling. This model improves client responsiveness

but adds complexity. The client must handle replies out of order. State management across request and reply is more complex. Error handling requires additional mechanisms.

## 4.3  Stateless vs. Stateful Servers

Stateless servers treat each request independently, with no memory of previous interactions. All information needed to process a request is contained in the request itself. Benefits include simpler server implementation, easier horizontal scaling since any server can handle any request, improved fault tolerance since server restart does not lose session state, and better load balancing since requests can be routed to any server.

Stateful servers maintain session state across requests from a given client. Session information is stored on the server and associated with client identity. Benefits include reduced request size since context need not be repeated, ability to implement multi-step transactions, support for server-push notification, and potentially improved performance through caching.

Many systems use a hybrid approach: stateless application servers with state stored in a separate data tier.

## 4.4  Connection Management

Connections between clients and servers may be managed in various ways. Per-request connections establish a new connection for each request, which is simple but has high overhead for frequent requests. Persistent connections maintain connections across multiple requests, reducing overhead but requiring connection management. Connection pooling maintains a pool of reusable connections, amortizing connection cost across requests. Multiplexed connections carry multiple concurrent requests on a single connection, as in HTTP/2.

## 4.5  Concurrency Models

Servers must handle multiple concurrent clients. Process-per-request creates a new process for each request, providing isolation but with high overhead. Thread-per-request creates a new thread for each request, which is lighter weight but requires careful synchronization. Thread pool uses a fixed pool of worker threads, bounding resource usage. Event-driven uses non-blocking I/O with event loops, supporting high concurrency with fewer threads. Hybrid approaches combine these models for different workload characteristics.

# 5  Constraints

The client–server style imposes constraints that define valid architectural configurations.

## 5.1  Connector Constraints

Clients are connected to servers through request/reply connectors. Direct client-to-client communication is not part of the base style. All client-server interaction flows through explicit connectors. This constraint ensures that communication is explicit and can be analyzed.

## 5.2 Role Asymmetry Constraint

Clients and servers have asymmetric roles. Clients initiate requests; servers respond to them. Servers do not initiate communication with clients in the basic model (though variations exist for server push). This asymmetry simplifies reasoning about system behavior.

## 5.3 Servers as Clients Constraint

Server components can be clients to other servers. This enables tiered architectures where each tier serves the tier above and consumes services from the tier below. A web server acts as a client to an application server; an application server acts as a client to a database server.

## 5.4 Attachment Constraints

Specializations may impose restrictions on attachments. Maximum clients per server may be limited by server capacity. Minimum servers per client may be required for availability. Exclusive attachment may require dedicated server instances for certain clients.

## 5.5 Tier Constraints

When components are arranged in tiers, additional constraints apply. Communication direction may be restricted so that requests flow up through tiers and replies flow down. Skip-tier communication may be prohibited, requiring each tier to interact only with adjacent tiers. Tier responsibilities may be constrained, such as limiting database access to the data tier.

## 5.6 Protocol Constraints

Clients and servers must agree on communication protocols. Data format compatibility requires shared understanding of request and response structure. Version compatibility requires managing protocol evolution. Security protocol requirements may mandate encryption or authentication.

# 6 What the Style is For

The client–server style supports several important architectural goals.

## 6.1 Separation of Concerns

The style promotes clean separation between service consumers and providers. Clients focus on user interaction and presentation. Servers focus on business logic and data management. This separation enables independent evolution, specialized optimization, and team organization aligned with concerns.

## 6.2 Modifiability and Reuse

By factoring out common services, the style promotes modifiability and reuse. Services can be modified without changing clients, as long as interfaces remain compatible. Services can be reused by multiple clients. New clients can be added without modifying servers. Implementation technologies can change behind stable interfaces.

## 6.3 Scalability

The style supports various scalability strategies. Horizontal scaling adds server instances to handle more clients. Vertical scaling adds resources to existing servers. Load balancing distributes requests across server instances. Caching reduces load on backend servers. The clear client-server boundary enables these strategies.

## 6.4 Availability

Server replication improves availability. Multiple server instances can handle requests if one fails. Failover mechanisms redirect clients to healthy servers. Load balancers can detect and route around failures. Stateless servers simplify replication.

## 6.5 Centralized Control

The style enables centralized control over data and business logic. Data integrity can be enforced at the server. Business rules are implemented consistently in one place. Access control is enforced centrally. Auditing and logging are simplified.

## 6.6 Security Analysis

The clear client-server boundary facilitates security analysis. Trust boundaries are explicit. Authentication can be enforced at server entry points. Authorization can be centralized in servers. Encryption can protect client-server communication. Attack surface is well-defined.

## 6.7 Performance Analysis

The style enables structured performance analysis. Server capacity can be measured and planned. Client load patterns can be characterized. Network bandwidth and latency can be analyzed. Bottlenecks can be identified at specific tiers.

## 6.8 Heterogeneity

Clients and servers can use different technologies. Web browsers, mobile apps, and command-line tools can access the same server. Servers can be implemented in different languages. This heterogeneity is enabled by well-defined interfaces.

# 7 Notations

Client–server architectures can be represented using various notations.

## 7.1 Box-and-Line Diagrams

Informal diagrams show clients and servers as boxes with lines representing connectors. Conventions include using distinct shapes for clients and servers, showing connector direction with arrows, grouping elements by tier, and labeling with component names and protocols.

These diagrams are widely understood but lack precise semantics.

## 7.2  UML Component Diagrams

UML provides formal notation for component architectures. Components represent clients and servers. Interfaces (lollipops) represent provided services. Required interfaces (sockets) represent service dependencies. Dependencies show client-to-server relationships. Deployment diagrams can show allocation to nodes.

## 7.3  UML Sequence Diagrams

Sequence diagrams show the temporal ordering of client-server interactions. Lifelines represent clients, servers, and other participants. Messages show requests and replies. Activation bars show processing time. Combined fragments show conditional and iterative behavior.

## 7.4  Architecture Description Languages

Formal ADLs provide precise semantics for client-server architectures. AADL supports embedded systems with client-server patterns. Wright provides formal specification of architectural styles. Rapide supports executable architecture specifications.

## 7.5  API Specifications

Modern systems often document client-server interfaces through API specifications. OpenAPI (Swagger) describes REST APIs. GraphQL schemas describe GraphQL APIs. Protocol Buffers describe gRPC services. WSDL describes SOAP web services.

These specifications serve as contracts between clients and servers.

## 7.6  Infrastructure Diagrams

Cloud and infrastructure diagrams show physical or virtual deployment. AWS, Azure, and GCP architecture diagrams use provider-specific iconography. Network diagrams show connectivity and security boundaries. Container orchestration diagrams show service deployment.

# 8  Quality Attributes

Client–server architectural decisions significantly affect system quality attributes.

## 8.1  Performance

Performance in client–server systems depends on multiple factors. Network latency dominates for remote communication; minimizing round trips is crucial. Server throughput limits the rate of request processing. Client responsiveness depends on asynchronous patterns and local processing. Caching at various levels reduces load and latency.

Performance tactics include connection pooling to reduce connection overhead, request batching to amortize latency across multiple operations, compression to reduce bandwidth usage, caching to avoid repeated processing, and asynchronous communication to improve client responsiveness.

## 8.2   Scalability

Scalability addresses growth in clients, data, or request rates. Horizontal scaling adds server instances, requiring stateless design or distributed state. Vertical scaling adds resources to existing servers, limited by single-machine capacity. Database scaling may require sharding or replication. Caching tiers offload read-heavy workloads.

Scalability is constrained by stateful components, shared resources, and synchronization requirements.

## 8.3   Availability

Availability ensures the system is operational when needed. Server replication provides redundancy. Load balancers detect and route around failures. Circuit breakers prevent cascade failures. Graceful degradation maintains partial functionality during failures.

Availability analysis considers failure modes, mean time between failures, mean time to recovery, and recovery procedures.

## 8.4   Security

Security in client–server systems addresses multiple concerns. Authentication verifies client identity using passwords, tokens, certificates, or multi-factor mechanisms. Authorization controls what authenticated clients can access. Data protection secures data in transit through TLS and at rest through encryption. Input validation prevents injection attacks. Audit logging records security-relevant events.

The client-server boundary is a natural point for security controls but also an attack surface.

## 8.5   Modifiability

Modifiability addresses the ease of changing the system. Interface stability enables independent client and server evolution. Versioning strategies manage interface changes. Backward compatibility protects existing clients. Loose coupling reduces change propagation.

Well-designed client-server interfaces enhance modifiability; tightly coupled interfaces impede it.

## 8.6   Reliability

Reliability addresses correct operation under various conditions. Error handling manages expected error conditions. Fault tolerance handles unexpected failures. Transaction support ensures data consistency. Idempotency enables safe retry of failed requests.

## 8.7   Portability

Portability addresses deployment across different environments. Standard protocols like HTTP enable diverse client platforms. Platform-independent data formats like JSON enable interoperability. Containerization enables consistent server deployment.

## 8.8 Testability

Testability addresses the ease of verifying system behavior. Clear interfaces enable unit testing of clients and servers independently. Mock servers support client testing. Service virtualization simulates complex backends. Contract testing verifies interface compatibility.

# 9 Common Client–Server Patterns

Several recurring patterns address common client–server challenges.

## 9.1 Two-Tier Architecture

In two-tier architecture, clients connect directly to database servers. Clients contain both presentation and business logic. The database server provides data storage and retrieval. This pattern is simple but has limitations. Business logic is distributed across clients, making updates difficult. Clients require database connectivity. Security enforcement relies on database access control.

This pattern is appropriate for simple applications with few clients and limited business logic.

## 9.2 Three-Tier Architecture

Three-tier architecture introduces an application tier between clients and data. The presentation tier handles user interface in clients. The application tier contains business logic in application servers. The data tier manages persistence in database servers.

Benefits include centralized business logic, thin clients, independent tier scaling, and clear security boundaries. This pattern is the foundation of most enterprise web applications.

## 9.3 N-Tier Architecture

N-tier architectures add tiers for specific concerns. Web tier handles HTTP and presentation. API tier exposes programmatic interfaces. Service tier implements business operations. Integration tier connects to external systems. Data tier manages persistence.

Additional tiers add separation but also complexity and latency.

## 9.4 Microservices Architecture

Microservices decompose the server side into independently deployable services. Each service owns its data and logic for a specific capability. Services communicate through APIs, often HTTP/REST or messaging. Services can be developed, deployed, and scaled independently.

This pattern enables organizational scaling and technology diversity but adds distributed systems complexity.

## 9.5 Backend for Frontend (BFF)

The BFF pattern creates dedicated backend services for different client types. A mobile BFF serves mobile clients. A web BFF serves browser clients. Each BFF tailors responses to its client's needs.

This pattern optimizes for client-specific requirements but adds backend components to maintain.

## 9.6 API Gateway Pattern

An API gateway provides a single entry point for clients. The gateway routes requests to appropriate backend services. It can handle cross-cutting concerns like authentication, rate limiting, and logging. It can aggregate responses from multiple services.

This pattern simplifies client interaction with complex backends but introduces a potential bottleneck.

## 9.7 Service Mesh Pattern

A service mesh handles service-to-service communication. Sidecar proxies handle communication, security, and observability. A control plane manages proxy configuration. The application code is freed from communication concerns.

This pattern adds infrastructure complexity but provides powerful operational capabilities.

## 9.8 CQRS Pattern

Command Query Responsibility Segregation separates read and write models. Commands modify state through one path. Queries read state through another path, often optimized differently. This pattern can improve performance and scalability for read-heavy workloads.

## 9.9 Event Sourcing with Client–Server

Event sourcing stores state changes as events. Clients send commands that generate events. Event stores persist the event stream. Read models project events for queries. This pattern provides audit trails and temporal queries but adds complexity.

# 10 Protocol Considerations

The choice of communication protocol significantly affects client–server systems.

## 10.1 HTTP and REST

HTTP is the dominant protocol for web-based client–server communication. REST architectural style uses HTTP methods semantically: GET for retrieval, POST for creation, PUT for update, and DELETE for removal. Resources are identified by URLs. Representations use formats like JSON or XML.

Benefits include universal client support, extensive tooling, cacheability, and statelessness. Limitations include verbosity, lack of streaming support in HTTP/1.1, and impedance mismatch with non-resource-oriented domains.

## 10.2 GraphQL

GraphQL provides a query language for APIs. Clients specify exactly what data they need. A single endpoint handles all queries. The schema defines available types and operations.

Benefits include client-specified responses reducing over-fetching, strong typing, and introspection. Limitations include complexity, potential for expensive queries, and caching challenges.

## 10.3   gRPC

gRPC uses Protocol Buffers for efficient binary serialization. It provides strongly-typed service definitions. Streaming is natively supported. HTTP/2 provides multiplexing and header compression.

Benefits include efficiency, streaming support, and code generation. Limitations include less tooling than REST, binary format complexity for debugging, and browser support challenges.

## 10.4   WebSocket

WebSocket provides full-duplex communication over a single connection. After an HTTP upgrade handshake, bidirectional messages can flow. This enables server push and real-time updates.

Benefits include low latency for frequent messages and server-initiated communication. Limitations include connection state management and scaling challenges.

## 10.5   Message Queues

Message queue protocols like AMQP and MQTT support asynchronous request/reply. Requests and replies are correlated through message identifiers. Queues provide buffering and delivery guarantees.

Benefits include decoupling, reliability, and support for offline clients. Limitations include added infrastructure and eventual consistency.

# 11   Error Handling

Robust error handling is essential for client–server systems.

## 11.1   Error Categories

Errors in client–server systems fall into several categories. Client errors result from invalid requests, authentication failures, or authorization denials. Server errors result from internal server failures, resource exhaustion, or downstream service failures. Communication errors result from network failures, timeouts, or protocol errors. Business errors result from domain rule violations distinct from technical errors.

## 11.2   Error Representation

Errors must be communicated from server to client. HTTP status codes provide standard error categories. Error response bodies provide details such as error codes, messages, and remediation hints. Structured error formats like RFC 7807 Problem Details standardize error representation.

## 11.3   Client Error Handling

Clients must handle errors appropriately. Retry logic handles transient failures with appropriate backoff. Circuit breakers prevent repeated calls to failing services. Fallback behavior provides degraded functionality when services are unavailable. User feedback communicates errors appropriately to users.

### 11.4 Server Error Handling

Servers must handle errors without compromising security or stability. Graceful degradation maintains partial functionality during partial failures. Error logging records errors for debugging and monitoring. Error abstraction hides internal details from clients to prevent information leakage. Transaction rollback maintains data consistency after failures.

### 11.5 Timeout Handling

Timeouts are critical for system stability. Client timeouts prevent indefinite waiting. Server timeouts prevent resource exhaustion. Cascading timeouts must be coordinated across tiers. Timeout values must balance responsiveness and false positives.

## 12 Security Considerations

Security is paramount in client–server architectures.

### 12.1 Authentication

Authentication verifies client identity. Username and password is the traditional approach with known vulnerabilities. API keys provide simple service-to-service authentication. OAuth and OpenID Connect provide delegated authorization and identity federation. Mutual TLS provides certificate-based authentication.

### 12.2 Authorization

Authorization controls access to resources and operations. Role-based access control assigns permissions to roles. Attribute-based access control makes decisions based on attributes of subject, resource, and environment. Policy-based systems externalize authorization logic.

### 12.3 Transport Security

Transport security protects data in transit. TLS encrypts communication and verifies server identity. Certificate pinning prevents man-in-the-middle attacks. Perfect forward secrecy protects past sessions if keys are compromised.

### 12.4 Input Validation

Input validation prevents injection attacks. All client input must be validated on the server. Parameterized queries prevent SQL injection. Output encoding prevents cross-site scripting. Schema validation ensures message structure.

### 12.5 Rate Limiting

Rate limiting prevents abuse and ensures fair resource allocation. Request rate limits cap requests per client per time period. Quotas limit total resource consumption. Throttling degrades service rather than failing completely.

## 12.6   Security Monitoring

Monitoring detects and responds to security incidents. Audit logging records security-relevant events. Intrusion detection identifies attack patterns. Alerting notifies operators of security events. Forensics support enables incident investigation.

# 13   Deployment Considerations

Client–server architecture affects deployment decisions.

## 13.1   Load Balancing

Load balancers distribute requests across server instances. Round-robin distributes requests evenly. Least-connections sends requests to least busy servers. Weighted distribution accounts for server capacity differences. Session affinity routes requests from a client to the same server.

Health checks detect failed servers. Active checks periodically probe servers. Passive checks monitor request failures. Failed servers are removed from the pool.

## 13.2   Server Replication

Server replication provides capacity and availability. Stateless servers are easily replicated; any instance can handle any request. Stateful servers require session replication or external session storage. Database replication requires consideration of consistency and conflict resolution.

## 13.3   Caching

Caching improves performance by storing frequently accessed data. Client-side caching stores data on the client, reducing requests. CDN caching stores static content at edge locations. Server-side caching stores computed results, reducing processing. Database caching stores query results, reducing database load.

Cache invalidation strategies ensure cached data remains current. Time-based expiration removes data after a set period. Event-based invalidation removes data when source data changes. Cache-aside patterns let the application manage the cache.

## 13.4   Geographic Distribution

Geographic distribution places servers near clients. Content delivery networks distribute static content. Regional deployments place application servers in multiple regions. Data replication strategies handle data distribution and consistency.

## 13.5   Container Orchestration

Container orchestration automates deployment and scaling. Kubernetes manages container lifecycle, scaling, and networking. Service discovery enables dynamic server location. Rolling updates enable zero-downtime deployments.

# 14    Examples

Concrete examples illustrate client–server concepts.

## 14.1    Web Application

A typical web application exhibits multiple client–server relationships. Browser clients request pages from web servers. Web servers request data from application servers. Application servers request data from database servers. JavaScript in browsers makes AJAX requests to API servers.

Considerations include session management (cookies, tokens), content delivery (static assets via CDN), API design (REST, GraphQL), and security (HTTPS, CORS, CSRF protection).

## 14.2    Mobile Application

Mobile applications present specific client–server challenges. Intermittent connectivity requires offline support and synchronization. Limited bandwidth requires efficient protocols and data formats. Diverse clients across iOS, Android, and various versions require API compatibility. Push notifications require server-to-client communication.

Considerations include API versioning, data synchronization, push notification infrastructure, and backend-for-frontend patterns.

## 14.3    Database Application

Database applications demonstrate classic two-tier client–server. Application clients connect to database servers using database protocols. Queries are requests; result sets are replies. Connection pooling manages database connections efficiently.

Considerations include query optimization, transaction management, connection management, and failover handling.

## 14.4    Microservices System

In microservices systems, services are both clients and servers. An API gateway serves external clients. Internal services call each other. Service discovery locates service instances. Circuit breakers handle service failures.

Considerations include service decomposition, inter-service communication (sync vs. async), distributed tracing, and data consistency.

## 14.5    Real-Time Application

Real-time applications require bidirectional communication. WebSocket connections enable server-initiated updates. Clients subscribe to topics or channels. Servers push updates to subscribed clients.

Considerations include connection scaling, message ordering, delivery guarantees, and reconnection handling.

# 15    Best Practices

Experience suggests several best practices for client–server systems.

## 15.1    Design for Failure

Assume that servers will fail and networks will be unreliable. Implement timeouts for all remote calls. Use circuit breakers to prevent cascade failures. Provide fallback behavior for unavailable services. Test failure scenarios regularly.

## 15.2    Prefer Stateless Servers

Stateless servers are easier to scale, deploy, and manage. Store session state externally in caches or databases. Include necessary context in requests. Use tokens rather than server-side sessions.

When state is necessary, manage it explicitly and plan for state loss.

## 15.3    Design Idempotent Operations

Idempotent operations produce the same result regardless of how many times they are executed. Idempotency enables safe retry of failed requests. Use unique request identifiers to detect duplicates. Design state changes to be idempotent where possible.

## 15.4    Version APIs Carefully

APIs evolve; versioning manages that evolution. Include version in API path or headers. Maintain backward compatibility within versions. Deprecate old versions with clear timelines. Document breaking changes clearly.

## 15.5    Implement Comprehensive Logging

Logging enables debugging, monitoring, and auditing. Log requests and responses with correlation identifiers. Include sufficient context for debugging. Structure logs for machine processing. Protect sensitive information in logs.

## 15.6    Monitor Everything

Monitoring provides visibility into system health. Track request rates, latencies, and error rates. Monitor resource utilization. Set up alerts for anomalies. Use distributed tracing for request flow visibility.

## 15.7    Secure by Default

Security should be built in, not added later. Use HTTPS everywhere. Authenticate all requests. Authorize at the server. Validate all input. Follow the principle of least privilege.

## 15.8    Document Interfaces

Clear interface documentation enables client development. Specify request and response formats. Document error conditions and codes. Provide examples. Keep documentation synchronized with implementation.

# 16    Common Challenges

Client–server systems present several common challenges.

## 16.1    Network Reliability

Networks are unreliable, introducing latency, packet loss, and partitions. Clients must handle timeouts and retries. Servers must handle duplicate requests. Systems must degrade gracefully during partitions.

## 16.2    Partial Failures

In distributed systems, partial failures are common. One server may fail while others continue. Downstream services may become unavailable. Clients must handle partial results and degraded functionality.

## 16.3    Consistency

Maintaining consistency across clients, servers, and data stores is challenging. Caches may serve stale data. Concurrent updates may conflict. Distributed transactions add complexity and reduce availability.

Strategies include accepting eventual consistency, using optimistic locking, implementing conflict resolution, and limiting the scope of transactions.

## 16.4    Latency

Network latency affects user experience and system throughput. Each client-server round trip adds delay. Deep call chains multiply latency. Geographic distribution increases baseline latency.

Strategies include reducing round trips through batching and aggregation, parallelizing independent requests, caching frequently accessed data, and placing servers near clients.

## 16.5    Scalability Bottlenecks

Various bottlenecks limit scalability. Single database instances limit write throughput. Stateful servers limit horizontal scaling. Shared resources create contention. Session affinity limits load distribution.

Strategies include sharding data, externalizing state, eliminating single points of contention, and designing for statelessness.

## 16.6    Version Compatibility

Clients and servers evolve independently, creating version compatibility challenges. Old clients must work with new servers. New clients should work with old servers during transitions. Breaking changes require coordinated rollout.

Strategies include versioned APIs, backward compatibility commitments, feature detection, and staged rollouts.

## 16.7   Debugging Distributed Systems

Debugging issues across clients and servers is challenging. Requests span multiple processes and machines. Logs are distributed. Timing-dependent bugs are hard to reproduce.

Strategies include correlation identifiers, centralized logging, distributed tracing, and comprehensive monitoring.

# 17   Conclusion

The client–server style provides a foundational pattern for distributed systems, enabling separation of concerns, scalability, and centralized control. By clearly defining clients that request services and servers that provide them, the style creates boundaries that enable independent development, deployment, and evolution.

Effective client–server architecture requires careful attention to many concerns: communication protocols, error handling, security, performance, and operational aspects. The patterns and practices described in this document provide guidance for navigating these concerns.

As systems grow in complexity, the basic client–server style is often composed with other styles and patterns. Microservices architectures decompose servers into smaller, independently deployable units. Event-driven patterns complement request/reply with asynchronous communication. The fundamental client–server concepts remain relevant even as the specific implementations evolve.

Investment in understanding client–server architecture pays dividends across virtually all modern software development, as the pattern underlies web applications, mobile backends, API platforms, and enterprise systems. Mastery of this style is essential for any practicing software architect.

# References

- Clements, P., Bachmann, F., Bass, L., Garlan, D., Ivers, J., Little, R., Merson, P., Nord, R., & Stafford, J. (2010). *Documenting Software Architectures: Views and Beyond* (2nd ed.). Addison-Wesley Professional.

- Bass, L., Clements, P., & Kazman, R. (2021). *Software Architecture in Practice* (4th ed.). Addison-Wesley Professional.

- Newman, S. (2021). *Building Microservices: Designing Fine-Grained Systems* (2nd ed.). O'Reilly Media.

- Richardson, C. (2018). *Microservices Patterns*. Manning Publications.

- Kleppmann, M. (2017). *Designing Data-Intensive Applications*. O'Reilly Media.

- Nygard, M. (2018). *Release It! Design and Deploy Production-Ready Software* (2nd ed.). Pragmatic Bookshelf.

- Fielding, R. T. (2000). *Architectural Styles and the Design of Network-based Software Architectures*. Doctoral dissertation, University of California, Irvine.