# Your Fun Guide to Cloudflare Durable Objects: Stateful Sorcery at the Edge!

Welcome, intrepid developer, to the fascinating world of Cloudflare Durable Objects! If you've ever wrestled with managing state in serverless applications, especially at the edge, then prepare to meet a technology that might just feel like a superpower. This guide is designed to be your fun, comprehensive companion to understanding what Durable Objects are, why they're a big deal, how they work their magic, and when you might want to unleash them in your own projects.

## 1. Durable Objects: Your New Superpower for Stateful Serverless

The digital landscape is increasingly shifting towards the edge, bringing applications and data closer to users for unparalleled speed and responsiveness. However, this shift presents its own set of challenges, particularly when it comes to managing state.

### Beyond Stateless: Why the Edge Needs State

Traditional serverless functions, including the foundational Cloudflare Workers, are celebrated for their stateless nature. Each invocation is typically independent, processing a request and returning a response without retaining memory of previous interactions. This model is fantastic for scalability and simplicity in many scenarios, like transforming requests, running A/B tests, or serving localized content.

However, many modern applications inherently require state – a memory of past events, user sessions, or shared data. When stateless edge functions need to access or modify state, they usually turn to external databases or caching layers. This introduces a potential wrinkle: if an application at the edge needs *strong consistency* (meaning it must always read the very latest data), all requests for that data might have to travel from the low-latency edge network to a more centralized data center.[1] This journey can reintroduce the latency that edge computing aims to minimize, creating a tension between speed and data accuracy. As one analysis puts it, "if we want strong consistency with our data store, all requests have to go from the low latency edge network to a data center somewhere, and we are back to not so great latency again. This is where Durables Objects come in".[1] This scenario perfectly illustrates the problem Durable Objects are designed to solve. The core promise of edge computing—low latency achieved by proximity to the user—can be undermined if state management requires frequent, long-distance data fetches. Durable Objects address this by bringing stateful capabilities directly into the edge compute layer.

### Meet Durable Objects: Compute + Storage = Awesome!

Enter Cloudflare Durable Objects (DOs). They are Cloudflare's innovative answer to the

challenge of stateful serverless computing at the edge. A Durable Object is a special kind of Cloudflare Worker that uniquely combines compute logic with its *own dedicated persistent storage*.[2] Think of it as a Worker that doesn't forget. Because this storage "lives together with the object, it is strongly consistent yet fast to access. Therefore, Durable Objects enable stateful serverless applications".[2]

Each instance of a Durable Object possesses a globally unique identity, ensuring that requests intended for a specific object always reach it. Crucially, each DO comes with its own private, transactional, and strongly consistent storage.[3] This tight coupling of compute and storage is what makes DOs a powerful primitive for building a new class of applications. They are available on both Free and Paid Cloudflare plans, and significantly, SQLite-backed Durable Objects are now accessible even on the Workers Free plan, lowering the barrier to entry for developers.[2]

A particularly apt description frames Durable Objects as "mini servers without the ops".[6] This analogy helps clarify their role: each DO is addressable by a unique ID, executes its code in a single-threaded fashion (per instance), maintains persistent state, and can coordinate multiple clients.[6] This server-like behavior, but without the traditional operational overhead of managing infrastructure, aligns perfectly with the serverless ethos. The "mini" aspect refers to their design philosophy: they are meant to "scale out, not up," meaning applications achieve scale by using many small, specialized DOs rather than a few large, monolithic ones. This model is ideal when an application needs an isolated, stateful entity to manage a specific piece of data or coordinate a particular set of interactions, much like one might use a dedicated microservice in other architectural paradigms.

## DOs vs. Regular Cloudflare Workers: What's the Big Deal?

To truly appreciate Durable Objects, it's helpful to see how they differ from standard Cloudflare Workers. While both run on Cloudflare's global edge network, their capabilities and ideal use cases diverge significantly. The core distinction lies in state management and data persistence.[1]

Here's a comparison to highlight these differences:

**Table 1: Cloudflare Workers vs. Durable Objects at a Glance**

| Feature | Cloudflare Worker | Durable Object | Why it Matters |
|---|---|---|---|
| **State Management** | Stateless | Stateful | DOs can remember information across requests and users for a specific instance. |
| **Primary Storage** | External (e.g., KV, D1, R2) | Built-in, co-located (SQLite or KV) | DOs have their own private storage, reducing latency and complexity for instance-specific data.[1] |

| | | | |
|---|---|---|---|
| **Consistency Model** | Varies by external store (e.g., KV is eventual) | Strong for its own storage (serializable with transactions) [7] | DOs provide strong guarantees for their data, simplifying logic that depends on up-to-date state. |
| **Execution Model** | Concurrent for multiple requests | Single-threaded *per DO instance* [1] | Simplifies concurrency management *within* a DO; requests to the same instance are processed sequentially. |
| **Instance Affinity** | No inherent affinity (any Worker can handle) | Affinity by unique ID | Requests for a specific DO ID are always routed to the same logical instance, ensuring state continuity. [4] |
| **Primary Use Case** | Edge logic, request/response transformation, caching | Coordination, stateful application logic, real-time interactions, session management | DOs are designed for tasks that require persistent state and coordination, while Workers excel at stateless edge computations. [2] |

This comparison underscores that Durable Objects are not merely "Workers with a database attached." They represent a distinct architectural pattern for managing state at the edge, offering strong consistency and low-latency access for individual object instances, a combination that is hard to achieve with traditional stateless functions relying on external databases.

# 2. The "Aha!" Moment: Unpacking the Magic of Durable Objects

Durable Objects possess a few core characteristics that, when understood, reveal why they are such a powerful tool for developers. These aren't just features; they are enablers of new application architectures and simplifiers of complex problems.

## Strong Consistency at the Edge: A Rare Gem

In the world of distributed systems, "strong consistency" is a highly desirable property. It means that any read operation will always return the most recently written data. Achieving this, especially across a globally distributed network like the edge, without a significant performance penalty, is a considerable challenge. Durable Objects tackle this head-on by

providing strongly consistent storage for each object instance.[2] The transactional storage API ensures that operations on a Durable Object's storage are atomic and isolated, guaranteeing that reads reflect the latest committed writes for that specific object.[7]

This strong consistency within a DO instance dramatically simplifies application logic. Developers don't need to contend with the complexities of eventual consistency—where data might be temporarily stale—or implement intricate state reconciliation mechanisms for data managed by that particular DO. When combined with the single-threaded execution model for each DO (which we'll explore next), it creates a highly predictable environment for state mutations. This reduces the cognitive burden on developers and minimizes a class of bugs often associated with concurrent state management. This makes DOs particularly well-suited for use cases where data integrity and predictable state changes are paramount, such as counters, leaderboards, booking systems, or managing the state of a collaborative session.

## Compute Meets Data: The "Zero-Latency" SQLite Advantage

One of the most exciting developments in the Durable Objects ecosystem is the deep integration of SQLite. When a Durable Object uses SQLite for its storage, the SQLite database engine runs *directly inside the Durable Object's process*, in the same thread as the application code.[8] This co-location is a game-changer. As described in Cloudflare's own materials, "Your application code runs exactly where the data is stored... your storage lives in the same thread as the application, requiring not even a context switch to access. With proper use of caching, storage latency is essentially zero".[8] This is a stark contrast to traditional serverless functions that access external databases over a network, where even a fast network introduces measurable latency.[1]

This "zero-latency" access isn't just a marginal performance improvement; it fundamentally alters the economics of data access within a Durable Object. Operations that would be prohibitively slow or complex if they involved network round-trips to a separate database become perfectly feasible. For instance, the common "N+1 selects" problem, where an application makes one query to fetch a list of items and then N subsequent queries to fetch details for each item, becomes much less of a performance concern when each query has near-zero latency.[8] This allows developers to write more straightforward, imperative-style data handling code within the DO, without needing to aggressively optimize for minimizing network calls to an external database.

SQLite-backed Durable Objects are now generally available and are the recommended option for new Durable Object classes.[2] They offer the power of relational data modeling, the familiarity of SQL querying, and robust features like Point-in-Time Recovery (PITR), which allows restoring an object's database to its state at any point in the last 30 days.[8] This combination of familiar, powerful database capabilities with the unique execution model of DOs makes them an incredibly attractive option for a wide range of stateful logic.

## The Beauty of Single-Threaded Concurrency (Yes, Really!)

It might sound counterintuitive in a world obsessed with parallelism, but each *individual* Durable Object instance operates on a single thread.[1] This means that requests sent to the

*same object instance* are processed sequentially, one after the other. While this might seem like a limitation, it's a deliberate design choice that brings profound benefits for developers. By ensuring single-threaded execution within an object, Durable Objects eliminate the need for complex locking mechanisms or manual management of race conditions when accessing or modifying that object's state.[3] The runtime handles the queuing and sequential processing of messages (requests) for each instance. This greatly simplifies the development of correct and robust stateful logic. As Cloudflare states, this model contributes to safety and correctness, allowing developers to "Choose three" – Easy, Fast, Correct – when it comes to managing state within a DO.[3]

It's crucial to understand that this single-threaded nature applies to an individual DO instance. Overall application scalability is not bottlenecked by this; rather, it's achieved by creating *many* Durable Object instances.[8] This is the "scale out, not up" philosophy in action. If an application needs to handle a high volume of concurrent operations on what is logically the same piece of data (e.g., a global vote counter), the architectural solution involves sharding that data or workload across multiple DOs, which can then be coordinated if necessary.[8] The unit of concurrency in the DO world is the object instance itself, not threads within an instance. This design aligns closely with the Actor model, where each actor processes messages sequentially, ensuring the integrity of its internal state.

# 3. Under the Hood: How Durable Objects *Actually* Work (The Fun Version!)

Now that we've seen *what* makes Durable Objects special, let's peek behind the curtain to understand *how* they operate. Their lifecycle, addressing mechanism, and communication patterns are key to their unique capabilities.

## A Day in the Life of a Durable Object (Creation, Activity, Hibernation, Eviction)

The lifecycle of a Durable Object is managed by the Cloudflare platform, abstracting away many of the complexities of traditional server management.3
Here's a typical day:

1. **Implicit Creation:** A Durable Object isn't explicitly "created" by a developer in the traditional sense. Instead, an instance is *implicitly created on its first access*.[3] When a Worker (or another DO) attempts to communicate with a DO using its unique ID, and that DO isn't currently active, the Cloudflare runtime instantiates it.
2. **Activity:** Once active, the DO processes incoming requests sequentially. It can access its in-memory state and its persistent storage (SQLite or KV).[3] The DO remains alive as long as it's actively processing requests.
3. **Hibernation:** After a period of inactivity (typically several seconds of being idle), the Durable Object may *hibernate*.[3] During hibernation, the DO's compute instance is effectively paused to save resources. A critical point here is that **in-memory state**

**(JavaScript variables) is lost upon hibernation.**[3] The object's constructor will run again the next time it's accessed after hibernating.[12] Therefore, any state that needs to persist across hibernations *must* be saved to its durable storage (e.g., ctx.storage) before hibernation and reloaded upon reactivation.[10]

4. **WebSocket Hibernation (A Special Case):** A remarkable feature is WebSocket Hibernation. Even if a DO hibernates, any WebSocket connections it was managing can remain open and connected from the client's perspective.[2] The DO is automatically woken up when a message arrives on one of these hibernated WebSockets. This is a huge boon for scalability and cost-efficiency in real-time applications.

5. **Eviction & Migration:** Over longer periods of inactivity, or due to platform maintenance, a hibernated DO might be evicted from the specific server it was running on. Its persistent storage, however, remains safe. Cloudflare also transparently manages the migration of DOs between healthy servers as needed, ensuring their continued availability without developer intervention.[3]

This lifecycle makes it clear that while DOs provide statefulness, they are not like traditional always-on servers. Hibernation is a core part of their efficiency model. Developers must design their DOs to correctly initialize their state from durable storage (e.g., in the constructor or on first access using a pattern like blockConcurrencyWhile) and persist any changes back to durable storage if they need to survive hibernation.[12] In-memory caching is for optimizing access to "hot" data during an active period, not for long-term persistence without explicit writes to storage.

## IDs, Stubs, and Bindings: Your Keys to the DO Kingdom

Communicating with a specific Durable Object requires a way to address it and a mechanism to send it messages. This is where IDs, stubs, and bindings come into play.

- **Unique Identifiers (IDs):** Every Durable Object instance is identified by a DurableObjectId. This ID is globally unique and serves as the permanent address for that specific object.[3] There are two primary ways to generate these IDs:
  - idFromName(name: string): This method, called on a DO namespace binding, generates a unique ID derived from the provided string. Given the same name (and called on the same class), it will always produce the same DurableObjectId.[16] This is useful when you need a predictable way to address an object, for example, a DO representing a specific document named "mydoc".
  - newUniqueId(): This method generates a cryptographically random, globally unique ID. These IDs generally offer better performance on first use because the system doesn't need to perform a global check to ensure another instance with the same name wasn't coincidentally created elsewhere.[16] However, if you need to access this same DO again later, you must store this randomly generated ID (e.g., in KV, D1, or another DO).
- **Bindings:** Cloudflare Workers (and other DOs that need to call DOs) gain access to Durable Object classes through *bindings* configured in the wrangler.toml (or wrangler.jsonc) file.[15] A binding essentially gives your Worker a named reference to a

DurableObjectNamespace for a particular DO class.
- **Stubs:** Once a Worker has an ID for a Durable Object, it gets a *stub* for that DO by calling the get(id) method on the namespace binding.[15] A DurableObjectStub is a client-side JavaScript object that acts as a proxy to the actual Durable Object instance. The Worker uses this stub to send messages (RPC calls) to the DO.[15]

It's important to note that Durable Objects do not receive HTTP requests directly from the public internet. Instead, an incoming request from a client first hits a regular Cloudflare Worker. That Worker then uses a binding to get a stub for the appropriate Durable Object and forwards the request (or a new RPC call based on it) to the DO.[15] This ID/stub mechanism is fundamental. The ID is the permanent, logical address, while the stub is the live, programmatic connection to the (potentially newly activated) physical instance. This system allows Cloudflare to reliably route requests for a given ID to the correct instance, wherever it might be running or hibernating, providing the illusion of a perpetually existing, addressable object.

## DOs as Actors: Orchestrating Your Distributed Systems

The architecture and behavior of Durable Objects align very closely with the **Actor programming model**.[3] In the Actor model:
- An **actor** is a fundamental unit of computation.
- Each actor has its own **private state** that no other actor can directly access.
- Actors communicate with each other exclusively by sending and receiving **asynchronous messages**.
- Each actor has a **mailbox** to queue incoming messages and processes them one at a time (sequentially).

This maps beautifully to Durable Objects:
- Each DO instance is like an actor.
- Each DO has its private, strongly consistent storage and in-memory state.
- DOs communicate via RPC calls (messages) from Workers or other DOs.
- The single-threaded execution model ensures that each DO processes its incoming requests (messages) sequentially.

This Actor model paradigm simplifies the design of concurrent and distributed systems by providing a higher-level abstraction for managing state and interaction.[3] It helps avoid many common concurrency pitfalls associated with shared memory. Cloudflare documentation explicitly draws parallels to established actor systems like Akka and Microsoft Orleans.[3] One detailed comparison highlights that, like Orleans Virtual Actors, Durable Objects have a "perpetual existence" (they are logical entities that always exist virtually, addressable by their ID) and benefit from "automatic instantiation" (the runtime activates them when needed).[19]

Thinking in terms of actors provides a powerful mental framework for designing applications with Durable Objects, especially for complex coordination tasks. Each DO becomes a well-defined unit of state and behavior. For example, a chat room can be a DO (an actor), a collaborative document can be a DO (an actor), and a user session can be a DO (an actor). This makes DOs more than just "stateful Workers"; they are building blocks for a specific style

of distributed system architecture known for its resilience and effectiveness in managing concurrency.

### Teamwork: How Workers and DOs Communicate (RPCs Demystified)

As mentioned, Durable Objects don't typically handle raw internet traffic. A Cloudflare Worker usually acts as the front door. This Worker receives an HTTP request, determines which DO should handle it (often by deriving a DO ID from the request path or headers), gets a stub for that DO, and then communicates with the DO.

This communication happens via Remote Procedure Calls (RPCs). The Worker can:

1. Call the fetch() method on the DO stub, passing it a Request object. The DO's fetch() handler then processes this request, similar to how a Worker handles an HTTP request. The URL for this Request object can be a conceptual one (e.g., http://do/pathname) used for routing within the DO.[18]
2. Define custom methods on its DO class (e.g., incrementCounter(), getUserProfile()). The Worker can then call these methods directly on the stub (e.g., stub.incrementCounter()). This is often a cleaner way to implement specific RPC interactions.[15]

Frameworks like Hono can simplify this interaction by providing routing capabilities within the Worker to elegantly map incoming HTTP requests to specific RPC calls on the Durable Object.[15] For example, an incoming POST /counter/increment to the Worker could be routed by Hono to call an increment() method on a Counter Durable Object. The two-step process is: Worker receives HTTP fetch, Worker makes RPC invocation to DO, DO processes and returns to Worker, Worker sends HTTP response back to client.[15]

# 4. The Treasure Inside: Exploring Durable Object Storage

At the heart of a Durable Object's statefulness is its private, persistent storage. This isn't just any storage; it's designed for speed, consistency, and ease of use, with SQLite now taking center stage.

### SQLite on the Edge: Your DO's Personal, Speedy Database

The integration of SQLite directly within Durable Objects is arguably one of its most compelling features.[8] As highlighted before, this means each DO can have its own private, relational database that runs in the same process and thread as its compute logic, leading to "essentially zero" latency for data access.[8]

Key advantages of SQLite-backed DOs include:

- **Relational Data Modeling:** The ability to define tables, columns, relationships, and indexes using standard SQL.
- **Powerful SQL Queries:** Leverage the full expressiveness of SQL for complex data retrieval and manipulation.
- **Strong Consistency & Transactions:** Within the context of a single DO's single-threaded execution, operations against its SQLite database are transactional and

strongly consistent.
- **"Zero-Latency" Access:** The co-location of compute and storage drastically reduces data access times.[8]
- **Point-in-Time Recovery (PITR):** SQLite-backed DOs offer the ability to restore an object's embedded SQLite database to its state at any point in the past 30 days, providing a safety net against accidental data corruption.[8]
- **General Availability & Recommendation:** SQLite in Durable Objects is now Generally Available (GA) and is the recommended storage backend for all new Durable Object classes.[2] Each SQLite-backed DO can have up to 10GB of storage.[3]

To use SQLite storage, developers enable it for a DO class via a migration step in their wrangler.toml (or wrangler.jsonc) configuration file, typically by adding the class name to the new_sqlite_classes array within a migration block.[12]

This embedded SQLite isn't just about adding a database; it's a paradigm shift for edge state. It combines the performance benefits of an embedded database with the managed, serverless nature of the Cloudflare platform. The familiarity of SQL and features like PITR make it a robust and developer-friendly solution for many per-instance state requirements, reducing the need to always reach for external, network-attached databases.

## The Classic Key-Value Store

Before the advent of direct SQLite integration, Durable Objects primarily offered a key-value (KV) storage API.[8] This API allows for storing, retrieving, and deleting data as key-value pairs.
- Even with SQLite-backed DOs, the KV API methods (get(), put(), delete(), list()) are still available. Interestingly, when used with an SQLite-backed DO, these KV operations are actually implemented by storing the data in a hidden SQLite table within that DO's database.[3]
- For DO classes that were created before the SQLite backend was the default (or explicitly use the older KV backend), this remains their primary storage mechanism. Cloudflare has indicated that KV-backed DOs remain for backward compatibility, with a migration path to SQL storage for existing classes planned for the future.[9]
- The KV store has its own limits, such as maximum key size (e.g., 2KB) and value size (e.g., 128KB).[7]

While SQLite is now the star, the KV API provides a simpler interface for use cases that naturally fit a key-value model.

## Understanding the Storage API and Strong Consistency

Regardless of whether using the explicit SQL API or the KV API, interaction with a DO's persistent storage happens through methods on the ctx.storage object, which is available in the DO's constructor and its methods.[16] Key methods include:
- ctx.storage.get(key): Retrieves a value by its key (KV API).
- ctx.storage.put(key, value): Stores a key-value pair (KV API).
- ctx.storage.delete(key): Deletes a key-value pair (KV API).
- ctx.storage.deleteAll(): Removes all data from the DO's storage.[12]

- ctx.storage.sql.exec(query, params): Executes an SQL query (SQL API).[12]
- ctx.storage.setAlarm(time) and ctx.storage.deleteAlarm(): Manage alarms.

A cornerstone of the DO storage API is its guarantee of **strong consistency and transactional operations** for that object's private storage.[2] Each storage method call (like put() or a series of operations within transaction()) is implicitly wrapped in a transaction, ensuring atomicity and isolation from other storage operations within that DO.[21]

To further prevent concurrency bugs that could arise from JavaScript's event-driven nature, the Durable Objects runtime, by default, employs "input gates" and "output gates." These gates typically pause the delivery of other I/O events (like new incoming requests) to the Object while a storage operation is in progress.[21] This ensures that, for example, a read operation isn't immediately followed by a write from an interleaved request before the first operation fully completes, which could lead to inconsistent state. This default serialization of storage access contributes significantly to the ease of writing correct stateful code within a DO.

# 5. More Cool Tricks: Special Features of Durable Objects

Beyond their core compute and storage capabilities, Durable Objects come with a few special features that unlock even more powerful and efficient application patterns.

## In-Memory State: For When Speed is Everything

Durable Objects can hold state directly in JavaScript variables within the class instance (e.g., this.counter = 0;).[2] Accessing this in-memory state is extremely fast, as it avoids any I/O operations to the persistent storage layer.

However, there's a crucial caveat: **in-memory state is volatile**. It is lost whenever the Durable Object is evicted from memory or hibernates due to inactivity.[3] When the DO is reactivated, its constructor runs again, and any previous in-memory values are gone unless they were explicitly reloaded.

Despite its transient nature, in-memory state is highly useful for:

- **Caching frequently accessed data:** Load data from persistent storage into an in-memory variable on first access or in the constructor, then serve subsequent requests from memory for the duration of the DO's active session.[3]
- **Batching operations:** Accumulate data or events in memory before writing them to persistent storage in a single, larger operation.
- **Storing short-lived session data:** Information that's only relevant while a user is actively interacting with the DO.

A common and recommended pattern is to initialize instance variables from persistent storage within the DO's constructor (often using ctx.blockConcurrencyWhile() to ensure initialization completes before other requests are processed) or on the first request that needs the data.[12] Any modifications to this state that need to survive eviction must be written back to ctx.storage. It's also a best practice to avoid using global variables for state within a DO class,

as different instances of the same DO class might inadvertently share or overwrite global variables, leading to unexpected behavior; instance variables (e.g., this.value) ensure that each DO instance maintains its own separate in-memory data.[13]

## WebSocket Hibernation: Scaling Real-Time Without Breaking a Sweat

Durable Objects are exceptionally well-suited for acting as WebSocket servers, managing real-time, bidirectional communication with clients.[3] This is ideal for applications like chat, live updates, and multiplayer games.

A standout feature here is **WebSocket Hibernation**.[2] Traditionally, maintaining a large number of persistent WebSocket connections can be resource-intensive and costly, as each connection might keep a server process active. WebSocket Hibernation allows a Durable Object to hibernate (and thus not incur compute duration charges) even if it has active WebSocket connections.[24] The connections themselves remain open from the client's perspective, managed by the Cloudflare edge. When a message arrives on a hibernated WebSocket, the runtime automatically wakes up the corresponding Durable Object to process it.

Key aspects of WebSocket Hibernation include:

- **Cost Reduction:** Significantly lowers costs for applications with many potentially idle WebSocket connections, as billable duration charges are not incurred during periods of inactivity.[24]
- **Scalability:** Enables applications to manage a massive number of concurrent WebSocket clients efficiently.
- **state.acceptWebSocket(webSocket):** This method (part of the native DO WebSocket API) is used to enable hibernation for a WebSocket connection, as opposed to the standard webSocket.accept() which does not support hibernation.[24]
- **serializeAttachment(attachment):** Allows a small amount of data (up to 2048 bytes, any type supported by the structured clone algorithm) to be associated with a WebSocket connection and survive hibernation.[24] This can be useful for storing minimal session context without needing a full storage read/write cycle just for hibernation. If larger amounts of data need to persist, the Storage API should be used, and perhaps a key or identifier stored via serializeAttachment.

WebSocket Hibernation is a game-changer for building scalable and cost-effective real-time applications on Cloudflare. It directly addresses one of the major operational and financial challenges of large-scale WebSocket deployments: the overhead of maintaining numerous, often idle, persistent connections. By allowing the DO compute to "sleep" while the connection endpoint remains alive at the edge, it ensures resources are consumed only when there's actual activity.

## Durable Object Alarms: Your DO's Built-in Scheduler

Durable Objects come equipped with an **Alarms API**, which allows an instance to schedule itself to wake up and execute code at a specific time in the future.[2] This is like having a built-in, per-object cron job or scheduler.

Alarms are incredibly useful for a variety of scenarios:
- **Periodic Tasks:** A DO can set a recurring alarm to perform routine maintenance, data aggregation, or checks.
- **Timeouts:** Implement timeouts for operations or sessions (e.g., automatically close an inactive game room).
- **Delayed Actions:** Schedule a task to be performed after a certain delay (e.g., send a follow-up notification).
- **Managing Time-To-Live (TTL):** A DO can set an alarm for its own expiration, performing cleanup tasks and potentially deleting its own storage when the alarm fires.[2]
- **Reliable Queuing:** As demonstrated in the initial architecture of Cloudflare Queues, alarms can ensure that tasks (like delivering messages from a queue) are reliably executed, even in the presence of failures, by scheduling retries or wake-ups.[26]

The Alarms API (ctx.storage.setAlarm(scheduledTime) and ctx.storage.deleteAlarm()) empowers DOs to manage their own temporal logic without relying on external scheduling services for tasks that are intrinsically tied to the state and lifecycle of that specific object. When an alarm fires, the DO's alarm() handler method is invoked.

# 6. When to Call in the DOs (And When to Keep Them on the Bench)

Durable Objects are a powerful and versatile tool, but like any specialized instrument, they excel in certain situations and might be less suitable for others. Understanding their sweet spots and limitations is key to leveraging them effectively.

## Sweet Spots: Ideal Use Cases for Durable Objects

Durable Objects shine in scenarios that require strong consistency for a specific piece of state, coordination among multiple clients or processes, or long-lived stateful interactions.
- Real-Time Wonders: Chats, Live Editing, Multiplayer Games
  This is a classic domain for DOs. Their ability to manage shared state, coordinate multiple WebSocket connections, and ensure consistent updates makes them ideal for:
  - **Chat Applications:** Each chat room or direct message thread can be a DO instance. It holds the message history (in its storage), manages the list of connected participants (potentially using in-memory state for active users and WebSockets), and broadcasts new messages to all participants.[2]
  - **Collaborative Document Editing:** A DO can represent a single document. It processes edits from multiple users sequentially, ensuring that all collaborators see a consistent version of the document and that conflicting edits are handled gracefully.[2]
  - **Multiplayer Games:** Each game session or interactive object within a game can be a DO. It manages the game state (player positions, scores, environment status) and synchronizes it across all players connected via WebSockets.[2] The "Multiplayer Doom" and "Full Tilt" examples showcase this.[28]

- Scalable SaaS: The "Database-per-Tenant" Dream
  For multi-tenant Software-as-a-Service (SaaS) applications, DOs offer an elegant way to achieve data isolation and scalability:
  - Each tenant (customer, user, organization) can be assigned their own dedicated Durable Object instance.[4]
  - If using SQLite-backed DOs, this effectively gives each tenant their own private, isolated SQLite database.
  - This pattern simplifies access control (no complex row-level security needed within a shared database), makes it easier to manage tenant-specific data and configurations, and allows for scaling by simply adding more DOs as new tenants onboard.[6] A blog post by Boris Tane demonstrates this "database per user" pattern, highlighting benefits like true isolation and simplified access control.[6]
- Smarter AI: Giving Your Agents a Memory
  As AI agents become more sophisticated, the need for persistent memory and stateful coordination becomes critical. Durable Objects are well-positioned to serve this need:
  - A DO can act as the "memory" for an AI agent, storing conversation history, user preferences, learned information, and task state.[2]
  - They can coordinate multi-step tasks for an agent, ensuring that long-running processes maintain their state and can be resumed.
  - Cloudflare itself highlights this as a top use case: "Build AI agents. You can give your agents memory to store data, the ability to coordinate tasks, and the smarts to make real-time decisions".[27] The recent announcement making DOs available on the free tier was explicitly linked to facilitating AI agent development.[5]
- **Other Powerful Applications:**
  - **Rate Limiters:** A DO can track request counts for a specific user, IP, or API key, enforcing usage limits with strong consistency.[25]
  - **Voting/Polling Systems:** Each poll or item being voted on can be a DO, atomically incrementing counts and ensuring data integrity.[28]
  - **Job/Queue Management:** As seen in the v1 architecture of Cloudflare Queues, a DO can manage a queue of tasks, leveraging alarms for reliable processing and retries.[26]
  - **Session Management:** Each user session can be managed by a DO, storing session data securely and efficiently.[30]

## Knowing the Limits: When DOs Might Not Be the Answer

While powerful, Durable Objects are not a universal solution. Their unique architecture comes with specific design considerations and limitations.
- The "Scale Out, Not Up" Philosophy Revisited
  A single Durable Object instance is single-threaded and has inherent throughput limits (e.g., a soft limit of around 1,000 requests per second).4 It's designed to be one of many. If an application anticipates extremely high traffic to the exact same logical entity (e.g., a single global counter being updated by millions of users simultaneously), that

load will need to be sharded across multiple DO instances. These shards might then require additional logic for coordination or aggregation.8 As stated, "When using Durable Objects, it's important to remember that they are intended to scale out, not up. A single object is inherently limited in throughput... To handle more traffic, you create more objects".8 This means developers must sometimes think more explicitly about data partitioning and distributed system patterns compared to using a traditional monolithic database that scales vertically.

- Key Limitations to Consider
  It's crucial to be aware of the platform limits when designing applications with Durable Objects 22:

**Table 2: Durable Object Limitations Quick Reference**

| Feature/Limit | Detail (Examples) | Implication/Consideration |
|---|---|---|
| **Storage per SQLite DO** | 10 GB [3] | Sufficient for many per-instance needs, but not for storing massive datasets within a single object. Sharding or external storage (R2) for larger data. |
| **CPU per Request** | Default 30 seconds, configurable up to 5 minutes [3] | Long-running computations within a single request handler can block other requests to that DO. Offload heavy tasks if possible. |
| **Request Rate per DO Instance** | Soft limit ~1000 requests/second [22] | High-throughput needs for a single logical entity require sharding across multiple DOs. |
| **Memory per Instance** | 128 MB (similar to Workers) [1] | Limits the amount of in-memory state. Rely on durable storage for larger datasets. |
| **Single-Threaded Execution (per DO)** | Requests to the same DO are processed sequentially. | Simplifies state management *within* the DO but means one slow request can delay others for that specific instance. |
| **Data Location** | A DO instance runs in one Cloudflare data center at a time. [3] | While DOs can be placed near users, all requests for a specific ID go to that one active location. Very high global read traffic might see latency. |
| **KV Storage Key/Value Sizes** | Keys: 2KB, Values: 128KB (for | Data needs to fit within these |

| | KV-backed DOs or KV API) [22] | limits if using the KV API directly for larger chunks. |
|---|---|---|

Given these characteristics, Durable Objects might not be the best fit if:
*   The primary need is to **store very large binary files** (Cloudflare R2 is designed for this).
*   The application requires **massive, concurrent write throughput to a single, indivisible piece of state** without the possibility of sharding.
*   The task is **purely stateless request/response transformation** (standard Cloudflare Workers are more efficient for this).
*   The application heavily relies on **built-in pub/sub mechanisms or presence features** that DOs don't offer natively (though these can often be built *using* DOs or by integrating with other services like an MQTT broker or Ably, as noted in a comparison [31]).
*   **Global, low-latency reads of the *same* data are paramount from all locations simultaneously.** Since a DO instance is active in one location at a time, reads from distant locations will incur network latency to reach that specific instance.[32] Caching strategies with Workers might be needed in front of such DOs for read-heavy global scenarios.

⚙ Copy And SaveShareAsk Copilot

Understanding these limitations is not about finding fault with Durable Objects, but about recognizing their design parameters. They are a specialized tool. Success with DOs often comes from embracing the "many small, coordinated objects" model and thoughtfully partitioning state and workload.

# 7. Durable Objects in Action: Inspiring Stories & Fun Builds

The true measure of a technology lies in what it enables. Durable Objects have been used by Cloudflare itself and by the wider developer community to build some truly innovative and impressive applications. These stories not only showcase the capabilities of DOs but also provide inspiration for what's possible.

## Cloudflare's Super Slurper: Turbocharging Data Migrations

Cloudflare's Super Slurper is a service designed to migrate large amounts of data into Cloudflare R2 storage. In a significant re-architecture, Cloudflare rebuilt Super Slurper from the ground up using its own developer platform – Cloudflare Workers, Queues, and Durable Objects – achieving up to a 5x improvement in transfer speeds.[33]
In this new architecture, Durable Objects play a critical role in managing the state and coordination of migration jobs [33]:
*   **Per-Account/Job Tracking:** Each customer account undergoing a migration was given

a dedicated Durable Object. This DO stored vital details for each migration job, such as bucket names, user options, and the overall job state. This ensured that each migration was managed in an organized and isolated manner.

- **Batch Management:** To handle the transfer of potentially billions of objects, "Batch DOs" were introduced. These DOs were responsible for managing all the objects queued for transfer within a batch, storing their transfer state, object keys, and other metadata.
- **Sharding for Scale:** Migrations can involve massive datasets. To handle this scale and to work around the per-object storage limit of SQLite-backed DOs (which was 10GB at the time of the blog post, and remains so [9]), the Super Slurper team implemented a sharding strategy for the Batch DOs. This distributed the request load and storage across multiple DO instances, preventing bottlenecks.

The Super Slurper case study is a powerful internal testament to the capabilities of Durable Objects. It demonstrates their use not just as simple state stores, but as active, stateful orchestrators in a complex, large-scale distributed workflow, coordinating with other serverless components like Workers and Queues. This success story should inspire confidence in using DOs for demanding, state-intensive applications.

## Gaming at the Edge: The Multiplayer Doom & "Full Tilt" Saga

The real-time, low-latency nature of Durable Objects makes them a natural fit for gaming applications.

- **Multiplayer Doom on Cloudflare Workers:** This impressive demo showcases a WebAssembly port of the classic game Doom, enhanced with multiplayer support, all running on Cloudflare's global network using Workers, WebSockets, Pages, and, crucially, Durable Objects.[28] While specific architectural details from a dedicated write-up are sparse in the provided materials, it's clear that Durable Objects would be instrumental in managing the game session state for each multiplayer match, synchronizing player actions, positions, and game events across all participants connected via WebSockets. Each active game instance would likely be managed by a unique DO.
- **"Full Tilt" - Winner of a Game Jam:** Developer Guido Zuidhof won the innovative category in a major game jam by creating "Full Tilt," a game where players use their smartphone as a Wiimote-like motion controller for a game running in their laptop's web browser.[29] The "secret sauce" was the combination of Jamstack principles with Cloudflare Workers and Durable Objects to achieve seamless, low-latency communication.
  - **Architecture:** Each game session (a "room") was represented by a Durable Object. This DO acted as a mini-server, relaying sensor data from the phone (connected via WebSocket) to the laptop (also connected via WebSocket).[29]
  - **Room Code Coordination:** A central "Room Hub" Durable Object was used to manage short, user-friendly room codes and map them to the actual game room DO instances. This allowed players to easily join the same session.[29]

- ○ **Low Latency:** Guido highlighted that DOs allowed him to "instantly spin up a personal gaming server anywhere in the world, as close to that player as possible," which was key to the game's responsiveness.[29]

These gaming examples demonstrate how DOs can enable globally distributed, interactive experiences. For multiplayer games, co-locating the game session state and coordination logic near the users is vital for minimizing latency. Durable Objects provide this capability within a serverless, on-demand model, lowering the barrier to entry for developing and deploying scalable, real-time interactive applications.

## Building Communities: Wildebeest & ActivityPub

Wildebeest was an open-source project by Cloudflare (now archived) to create an ActivityPub and Mastodon-compatible server, allowing anyone to operate their own Fediverse identity and server on their domain with minimal infrastructure.[28] It was built using a suite of Cloudflare services: Workers, Pages, Durable Objects, Queues, and D1.

While detailed specifics of DO usage in Wildebeest are not exhaustively documented in the provided snippets, their role can be inferred from the project's nature and DO capabilities [36]:

- **Managing User State:** Each user account, its profile, list of followers/following, and notifications could be managed by or associated with Durable Objects.
- **Handling Posts and Interactions:** Individual posts ("toots") and their interactions (boosts, replies) might involve DOs for state management and real-time updates.
- **Caching:** The Wildebeest README mentions a switch from KV to DO for caching, indicating DOs were used to store frequently accessed data for performance.[36]

Wildebeest, even as an archived project, serves as an ambitious example of using DOs to build complex, stateful, server-side applications for decentralized social networking.

## Synchronized Experiences: TinyBase and Live Cursors

Durable Objects are also proving valuable for building collaborative tools that require synchronized state across multiple clients.

- **TinyBase Integration:** TinyBase, a library for structured local data, leverages Durable Objects to create full-stack, real-time, multi-device collaborative applications.[14] In this pattern, a Durable Object functions as a WebSocket server and also has its private, transactional, and strongly consistent storage. Client-side TinyBase Stores synchronize their data with the DO, which persists the canonical state. This provides a complete solution for building complex collaborative apps where data consistency and real-time updates are key.[14]
- **Live Cursors Tutorial:** A Cloudflare tutorial demonstrates building a "Live Cursors" feature (where multiple users see each other's mouse cursors in real-time on a shared page) using Next.js for the frontend, RPC for communication, and Durable Objects for the backend.[37] In this setup, the Durable Object acts as the central coordination hub. When a user moves their cursor, the frontend sends an update (via RPC through a Worker) to the DO. The DO then broadcasts this new cursor position to all other clients connected to the same collaborative session, likely via WebSockets managed by the

DO.[37]

## Other Cool Implementations (Briefly)

- **Cloudflare Queues (v1 Architecture):** In its initial version, each individual message queue in Cloudflare Queues was implemented as a single Durable Object.[26] This DO used its storage to hold messages and its Alarms API to schedule message delivery and retries. While Queues v2 evolved to use multiple sharded DOs per queue for better throughput and lower latency, the v1 architecture showcased DOs' capability for reliable, stateful task processing.[26]
- **NBA Finals Polling and Predictor:** This demo application uses Cloudflare Workers AI, Pages, Hono, and Durable Objects to create a stateful polling experience.[28] Durable Objects are used to keep track of user votes for different basketball teams and to store data for generating personalized predictions.
- **Seat Booking App Tutorial:** This tutorial guides developers in building a transactional seat booking application using SQLite-backed Durable Objects.[37] The DO manages the seat inventory, ensuring that bookings are processed atomically and preventing overbooking, leveraging the strong consistency and transactional capabilities of its embedded SQLite database.

These diverse examples highlight the versatility of Durable Objects, from powering internal Cloudflare services to enabling indie game development and sophisticated collaborative platforms.

# 8. Your Turn! Getting Started with Durable Objects

Feeling inspired? Getting your hands dirty with Durable Objects is the best way to truly understand their power. Here's a quick guide to get you started.

## A Quick Peek at Setting Up Your First DO

Setting up a Durable Object involves a few key steps, primarily using Cloudflare's command-line tool, Wrangler.

1. **Prerequisites & Plans:**
    - Historically, Durable Objects required a Workers Paid plan. However, a significant recent development is that **SQLite-backed Durable Objects are now available on the Workers Free plan**, along with generous allowances for requests, duration, and storage.[2] This dramatically lowers the barrier to entry for experimentation and building smaller projects. For KV-backed DOs or higher limits, the Workers Paid plan is still necessary.[39]
    - You'll need wrangler, the CLI for Cloudflare Workers, installed and configured. Version 2.0 or later is generally recommended.[39]
2. **Create a Worker Project:**
    - The easiest way to start is by using create-cloudflare. You can run a command like npm create cloudflare@latest my-durable-object-project --

--type=worker-durable-object or follow the interactive prompts, choosing a template that includes Durable Objects.[16] This will set up a basic project structure with a Worker script and a Durable Object class.

3. **Define Your Durable Object Class:**
   - A Durable Object is defined as a JavaScript or TypeScript class that typically extends DurableObject (though it's not strictly required to extend it, the DurableObjectState is passed to the constructor).[15]
   - The constructor receives ctx (or state in older examples), which provides access to the object's storage (ctx.storage), and env, which provides access to bindings (like KV namespaces or other DOs).
   - You'll implement methods on this class to handle logic. A common method is fetch(request) if you want the DO to handle HTTP-like requests from a Worker, but you can also define any other custom methods for RPC.

⚙ TypeScript

```typescript
// Example Durable Object class
export class Counter extends DurableObject {
  constructor(ctx: DurableObjectState, env: Env) {
    super(ctx, env);
    // Initialize state from storage if needed
    ctx.blockConcurrencyWhile(async () => {
      let storedValue = await this.ctx.storage.get<number>("value");
      this.value = storedValue |
```

Copy And SaveShareAsk Copilot

```typescript
| 0;
    });
  }
```

```typescript
  value: number = 0;

  async increment(amount: number = 1) {
    this.value += amount;
    await this.ctx.storage.put("value", this.value); // Persist state
    return this.value;
  }

  async getValue() {
    return this.value;
  }

  // Optional: if called via stub.fetch()
```

```
  async fetch(request: Request) {
    const url = new URL(request.url);
    if (url.pathname === "/increment") {
      const amount = parseInt(url.searchParams.get("amount") |
```

🔗 Copy And SaveShareAsk Copilot

```
| "1");
const newValue = await this.increment(amount);
return new Response(Counter incremented to: ${newValue});
} else if (url.pathname === "/value") {
const currentValue = await this.getValue();
return new Response(Current counter value: ${currentValue});
}
return new Response("Not found", { status: 404 });
}
}
```

4. **Configure Bindings and Migrations in wrangler.toml:**
   - Your wrangler.toml (or wrangler.jsonc) file needs to declare a binding for your Durable Object class. This allows your Worker to create/access instances of this DO class.[15]
   - You also need to define a migration if you're creating a new DO class, especially if it's an SQLite-backed DO.[12]

⚙ Ini, TOML
```
# wrangler.toml example
name = "my-durable-object-project"
main = "src/worker.ts" # Your Worker entrypoint
compatibility_date = "YYYY-MM-DD"

[[durable_objects.bindings]]
name = "COUNTER_DO" # How your Worker will refer to this DO namespace
class_name = "Counter" # The exported class name of your DO

# To enable SQLite storage for the "Counter" class
[[migrations]]
tag = "v1-counter-sqlite" # A unique tag for this migration
new_classes = ["Counter"] # Old way, for KV-backed
new_sqlite_classes = ["Counter"] # New way, for SQLite-backed DOs
```

Copy And SaveShareAsk Copilot*Note: For SQLite, ensure new_sqlite_classes is used as per current best practices.[2]*

5. **Write Your Worker Code to Interact with the DO:**
   - In your Worker's fetch handler (or other event handlers), you'll use the binding name to get a namespace object.

- From the namespace, get an ID for your DO (e.g., env.COUNTER_DO.idFromName("myCounter")).
- Get a stub for that ID (env.COUNTER_DO.get(id)).
- Call methods on the stub (e.g., await stub.increment(5) or await stub.fetch("http://do/value")).

⚙ TypeScript

```typescript
// Example Worker code (src/worker.ts)
export default {
  async fetch(request: Request, env: Env, ctx: ExecutionContext): Promise<Response> {
    const id = env.COUNTER_DO.idFromName("singletonCounter");
    const stub = env.COUNTER_DO.get(id);

    // Example: Call a custom RPC method
    // const newValue = await stub.increment(1);
    // return new Response(`Counter is now: ${newValue}`);

    // Example: Forward request to DO's fetch handler
    return stub.fetch(request);
  },
};

export interface Env {
  COUNTER_DO: DurableObjectNamespace<Counter>; // Matches binding name
}
```

Copy And SaveShareAsk Copilot

6. **Develop Locally and Deploy:**
   - Wrangler allows for local development and testing of Workers and Durable Objects.
   - When ready, deploy using wrangler deploy.

This provides a basic workflow. The Cloudflare documentation offers detailed tutorials and examples for various scenarios.[16]

# Tips for Worker-DO Interaction

- **RPC Flow:** Remember the typical flow: Client -> Worker -> DO -> Worker -> Client. The Worker acts as the intermediary.[15]
- **Hono for Routing:** For more complex Workers that interact with DOs, consider using a lightweight routing library like Hono. Hono can define clean API routes in your Worker that map to specific methods or fetch calls on your Durable Objects, making the Worker code more organized and maintainable.[15]
- **Validate in Worker:** To save on Durable Object invocation costs and processing time, perform as much request validation (authentication, schema validation, etc.) as possible in the Worker *before* calling the Durable Object.[24] Only valid and necessary requests should reach the DO.

- **Error Handling:** Implement robust error handling. Exceptions from a DO (either from user code or DO infrastructure) will propagate to the calling Worker. The Worker should catch these and respond appropriately to the client.[40] Some errors might be retryable.

# 9. Going Deeper: Advanced Durable Object Patterns

Once comfortable with the basics, developers can explore more advanced patterns to build highly sophisticated and scalable systems with Durable Objects.

## DO-to-DO Communication Patterns

Durable Objects are not limited to being called only by Workers; they can also communicate directly with other Durable Objects.[3] If a Durable Object class has bindings to other Durable Object namespaces defined in its wrangler.toml, it can obtain a stub for another DO instance and call its methods, just like a Worker would.

This capability unlocks several powerful architectural patterns:

- **Coordinator DOs:** A "manager" DO can orchestrate the work of several "worker" DOs. For example, a ProjectCoordinatorDO might manage a list of TaskDO instances, delegating operations to them and aggregating status.
- **Fan-out/Fan-in:** A common pattern for scaling operations. An initial DO might receive a request that involves processing a large dataset. It can then "fan out" by creating or calling multiple other DOs, each responsible for a subset of the data. These worker DOs process their parts in parallel. Later, their results can be "fanned in" – collected and aggregated by the original DO or another designated aggregator DO. This approach was conceptually mentioned for scaling a high-traffic vote counter.[8]
- **Chained DOs:** A request or workflow might pass through a sequence of DOs, each performing a specific transformation or stateful operation before passing it to the next.
- **Hierarchical Actor Systems:** Developers can design systems where DOs form parent-child relationships, with parent DOs supervising and delegating to child DOs.

The ability for DOs to directly message each other (via their stubs) allows for the creation of complex graphs or networks of collaborating actors. This is essential for implementing patterns like distributed consensus (among a small, controlled group of DOs), sharded counters with an aggregator DO, or intricate workflow orchestrations where different DOs handle distinct stages of a process. This elevates DOs from being merely isolated stateful entities to becoming components of potentially complex, self-organizing distributed systems, all running on the Cloudflare edge.

## Navigating Concurrency: The allowConcurrency Setting

By default, the Durable Objects runtime serializes access to an object's storage to prevent common concurrency bugs. When a DO is performing a storage operation (e.g., get(), put(), sql.exec()), the system typically pauses the delivery of new I/O events (like incoming requests or alarm triggers) to that specific DO instance until the storage operation completes.[21] This is managed by "input gates."

However, the storage API methods offer an allowConcurrency: true option. Setting this to true

opts out of this default serialization behavior, allowing other I/O events to be delivered to the DO and potentially begin execution *while a previous storage operation is still in progress.*[21]
**Implications and Risks of allowConcurrency: true:**
- **Potential for Race Conditions:** If one event handler starts a storage read, and another concurrent event handler (allowed because of allowConcurrency: true) modifies the same data before the first read completes, the first handler might operate on stale data. Similarly, read-modify-write sequences can become corrupted.[21]
- **Non-Atomic Multi-Step Operations:** If a logical operation involves multiple storage calls that must appear atomic, allowing concurrency can break this atomicity unless carefully managed by the developer.
- **Increased Code Complexity:** Reasoning about the state of the DO becomes significantly more complex, as the developer must account for all possible interleavings of operations.
- **When to Consider It (Cautiously):** This option might be considered for long-running, read-only storage operations where the DO needs to remain responsive to other, unrelated events. Or, if the developer is implementing their own sophisticated concurrency control mechanisms within the DO. However, it should be used with extreme caution and a deep understanding of the potential pitfalls.[21]

For most use cases, the default behavior (implicitly allowConcurrency: false) provides a safer and simpler programming model.

## Sharding Strategies for Massive Scale

The "scale out, not up" principle is fundamental to Durable Objects.[8] While a single DO can handle a respectable load, truly massive scale for a single logical entity (like a global counter accessed by millions, or an extremely popular chat room) often requires **sharding**. This means partitioning the data and/or workload across multiple DO instances.
Examples of sharding strategies:
- **Identifier-Based Sharding:** For a chat application, instead of one DO for "SuperPopularChatRoom", one might shard by user ID (e.g., users A-M go to SuperPopularChatRoom_Shard1, users N-Z go to SuperPopularChatRoom_Shard2). These shard DOs would then need a mechanism to exchange messages if users in different shards need to communicate.
- **Random/Round-Robin Sharding for Counters:** To implement a high-throughput global counter, increments could be distributed across N counter DOs. Each increment request goes to a randomly chosen shard DO. An aggregator (another DO or a Worker) can then periodically read the values from all shard DOs and sum them to get the global total.[8]
- **Cloudflare's Own Examples:**
  - **Super Slurper:** Used sharding for its "Batch DOs" to manage billions of objects and overcome the 10GB per-object storage limit.[33]
  - **Cloudflare Queues (v2):** The v2 architecture explicitly uses multiple "Storage Shard Durable Objects" per queue, distributed across regions, to improve latency

and throughput over the v1 single-DO-per-queue model. A Worker uses a "Shard Map" (stored in Workers KV) to load balance requests across available shards.[26] Sharding is the path to (almost) unlimited scale with DOs. While it introduces complexities (how to assign requests to shards, how to aggregate data across shards, potential for hot shards), it's the standard way to push beyond the limits of a single instance. The "database-per-user" pattern discussed earlier is a natural form of sharding, where the user ID is the shard key.[6] The platform provides the building blocks (many DOs), but the specific sharding logic is typically application-defined.

# 10. The Bottom Line: Understanding Durable Objects Pricing

Understanding how Durable Objects are priced is crucial for designing cost-effective applications. The model has several components, and recent changes have made DOs more accessible than ever.

## How Costs are Calculated

Previously, Durable Objects usage was primarily tied to the Workers Paid plan, which has a minimum charge (e.g., $5 USD per month per account) that includes a base allotment of Workers usage, KV operations, and Durable Objects usage.[41]

However, a major recent announcement (April 2025) states that **Durable Objects are now available on the Workers Free plan**.[2] This is particularly aimed at helping developers build AI agents and experiment with stateful applications without an initial financial commitment. The free tier includes substantial monthly allocations for [5]:

- **Durable Object Requests:** e.g., 100,000 per day.
- **Compute Duration:** e.g., 13,000 GB-seconds per day.
- **SQL Storage:** e.g., 5 GB total.

For usage beyond the free tier allowances or for features/limits exclusive to the Paid plan, the billing metrics for Durable Objects generally include [42]:

1. **Requests:**
   - This counts HTTP requests made to a DO's fetch() handler, RPC method calls, WebSocket messages, and Alarm invocations.
   - WebSocket messages have a special billing consideration: for example, 20 incoming WebSocket messages might be billed as 1 request for billing purposes (this ratio can vary).[42] Outgoing WebSocket messages are typically not charged as requests.
2. **Duration (GB-seconds):**
   - This measures how long Durable Objects are actively consuming compute resources. It's calculated based on the memory allocated to the DO (typically 128MB) multiplied by the time it's active.
   - A key factor here is **hibernation**. When a DO is hibernating (including WebSocket Hibernation), it does *not* accrue duration charges.[24] This makes features like

WebSocket Hibernation critical for cost optimization in real-time applications with many idle connections.

3. **Storage:**
   - **SQLite Storage Backend:**
     - **Rows Read:** Charged per million rows read from the SQLite database.
     - **Rows Written:** Charged per million rows written to the SQLite database.
     - **Stored Data:** Charged per GB-month for the amount of data stored in the SQLite databases across all DOs.
   - **Key-Value Storage Backend:**
     - **Read Request Units:** Based on 4KB units of data read.
     - **Write Request Units:** Based on 4KB units of data written. setAlarm is also billed as a write unit.
     - **Delete Requests:** Typically billed per deletion operation.
     - **Stored Data:** Charged per GB-month for data stored.

The pricing model, particularly the duration charges, incentivizes developers to design DOs that are active only when necessary. Efficient use of hibernation and careful consideration of in-memory versus persistent storage are key to managing costs. The new free tier significantly lowers the barrier for getting started and for running smaller applications.

## Example Scenarios to Demystify Billing

The Cloudflare documentation provides illustrative examples of how different usage patterns affect costs.[42] Let's consider the spirit of those:

- **Simple HTTP DO (Low Activity):** Imagine a single DO used for a coordination task, called 1.5 million times a month, active for a total of 1,000,000 seconds (about 277 hours) with 128MB memory.
  - Requests: (1.5M - 1M included on a paid plan) * $0.15/M = $0.075.
  - Duration: 1,000,000s * (128MB / 1024MB_per_GB) = 125,000 GB-s. If the paid plan includes 400,000 GB-s, this duration cost would be $0.00.
  - Total (excluding storage and base plan fee): ~$0.075.
  - *With the new free tier, if this usage fits within the daily/monthly free allowances, the cost could be $0.*
- **Many Active DOs (Moderate Activity):** Consider 100 DOs, each active for 8 hours a day for 30 days, handling a total of 3.75 million requests.
  - Requests: (3.75M - 1M included) * $0.15/M = $0.41.
  - Duration: 100 DOs * 8 hr/day * 3600 s/hr * 30 days * (128MB/1024MB_per_GB) = 10,800,000 GB-s.
  - Duration Cost: (10.8M GB-s - 0.4M GB-s included) * $12.50/M GB-s = $130.
  - Total (excluding storage and base plan fee): ~$130.41.
- **WebSocket Chat App (No Hibernation):** 100 DOs, each managing a WebSocket connection that sends/receives 1 message per second, 24/7 for 30 days.
  - WebSocket Connection Requests: 100.
  - WebSocket Messages: 1 msg/s * 100 conns * 3600 s/hr * 24 hr/day * 30 days = 259,200,000 messages.

- Billed Requests (20:1 ratio): 100 + (259.2M / 20) = ~12.96M requests.
- Request Cost: (12.96M - 1M included) * $0.15/M = ~$1.79.
- Duration (active 24/7): 100 DOs * 24 hr/day * 3600 s/hr * 30 days * (128MB/1024MB_per_GB) = 32,400,000 GB-s.
- Duration Cost: (32.4M GB-s - 0.4M GB-s included) * $12.50/M GB-s = ~$400.
- Total (excluding storage and base plan fee): ~$401.79.
- **WebSocket Chat App (With Hibernation):** Same as above, but assume hibernation means each DO is only truly active for 1 second per minute to process messages, then hibernates.
  - Request Cost: Same, ~$1.79.
  - Duration (active 1s per minute): 100 DOs * (1/60 active fraction) * 24 hr/day * 3600 s/hr * 30 days * (128MB/1024MB_per_GB) = 540,000 GB-s.
  - Duration Cost: (540,000 GB-s - 400,000 GB-s included) * $12.50/M GB-s = ~$1.75.
  - Total (excluding storage and base plan fee): ~$3.54.

These examples dramatically illustrate the cost-saving impact of WebSocket Hibernation. Storage costs (SQLite rows read/written, data stored) would be additional, depending on the application's data patterns. Always refer to the official Cloudflare pricing page for the most current details and free tier allowances.

# 11. Are Durable Objects Right for You? The Final Takeaway

Cloudflare Durable Objects are a genuinely innovative and powerful addition to the serverless landscape. They offer a unique solution to the persistent challenge of managing state at the edge, enabling new classes of applications that require strong consistency, real-time coordination, and low-latency stateful interactions.

## Recap: The Awesome, The Good, and The Things to Watch Out For

Let's distill our journey into a quick summary:
- **The Awesome:**
  - **Strongly Consistent State at the Edge:** A rare and valuable capability, simplifying development for state-dependent logic.[2]
  - **"Zero-Latency" SQLite Storage:** Co-located, relational database power inside each DO makes data access incredibly fast and familiar.[8]
  - **Simplified Concurrency (Per Object):** The single-threaded model per instance removes many common concurrency headaches for that object's state.[3]
  - **WebSocket Hibernation:** Massively scalable and cost-effective real-time communication by allowing idle connections without continuous compute cost.[2]
  - **Actor Model Paradigm:** Provides a robust mental model for designing distributed, stateful systems.[3]
  - **Now on the Free Tier:** Greatly increased accessibility for experimentation and

smaller projects.[5]
- **The Good:**
  - **Enables New Application Classes:** Ideal for real-time collaboration, multiplayer games, sophisticated SaaS multi-tenancy, and AI agent memory.[2]
  - **Built-in Scheduling (Alarms):** Useful for per-object timed tasks, TTLs, and reliable queuing patterns.[3]
  - **Managed Lifecycle:** Cloudflare handles instantiation, hibernation, and migration, reducing operational burden.[3]
- **Things to Watch Out For (Design Considerations):**
  - **Scale Out, Not Up:** A single DO has throughput limits; design for many small, coordinated objects if massive scale for one logical entity is needed.[8]
  - **Per-Object Limits:** Be mindful of storage (10GB/SQLite DO), CPU, memory, and request rate limits per instance.[22]
  - **Hibernation and State:** Understand that in-memory state is lost on hibernation; persistent state must be explicitly saved to and loaded from ctx.storage.[3]
  - **Potential Sharding Complexity:** If sharding is required for extreme scale, the logic for partitioning, routing, and aggregation is application-defined and can add complexity.
  - **Data Egress for a Single Instance:** A DO runs in one location at a time. If global clients all need to *write* to the same DO or read its very latest state, they'll all be routed to that one location, which might not be optimal for all users' latency.

## Making an Informed Decision for Your Next Project

Durable Objects are not a universal replacement for all stateful backends or traditional databases. They are a specialized tool that excels when the problem domain aligns with their strengths:
- Does the application require **strong consistency for specific, identifiable units of state** (like a user session, a document, a game room, a tenant's data)?
- Is **coordination among multiple clients or processes around a shared piece of state** a core requirement?
- Would **low-latency access to this state from edge compute** significantly benefit the application?
- Does the **Actor model** resonate as a good way to structure the application's components?

If the answers to these questions are "yes," then Durable Objects are certainly worth serious consideration. The recent expansion of their availability onto the Cloudflare Free plan makes it easier than ever to experiment, build proofs-of-concept, and even run production applications without upfront costs.[5]

They represent a significant step forward in making stateful serverless at the edge not just possible, but practical, powerful, and even fun to work with. Happy building!

**Works cited**

1. Cloudflare Durable Objects - Weird and Wonderful - Magnus Wahlstrand, accessed May 20, 2025, https://wahlstrand.dev/posts/cloudflare-durable-objects---weird-and-wonderful/
2. Overview · Cloudflare Durable Objects docs, accessed May 20, 2025, https://developers.cloudflare.com/durable-objects/
3. What are Durable Objects? - Cloudflare Docs, accessed May 20, 2025, https://developers.cloudflare.com/durable-objects/what-are-durable-objects/
4. Durable Objects (DO) — Unlimited single-threaded servers spread across the world, accessed May 20, 2025, https://www.lambrospetrou.com/articles/durable-objects-cloudflare/
5. Cloudflare unveils major AI Agent development tools - PPC Land, accessed May 20, 2025, https://ppc.land/cloudflare-unveils-major-ai-agent-development-tools/
6. One Database Per User with Cloudflare Durable Objects and Drizzle ORM | Boris Tane, accessed May 20, 2025, https://boristane.com/blog/durable-objects-database-per-user/
7. Choosing a data or storage product. - Workers - Cloudflare Docs, accessed May 20, 2025, https://developers.cloudflare.com/workers/platform/storage-options/
8. Zero-latency SQLite storage in every Durable Object - The Cloudflare Blog, accessed May 20, 2025, https://blog.cloudflare.com/sqlite-in-durable-objects/
9. SQLite in Durable Objects GA with 10GB storage per object - Cloudflare Docs, accessed May 20, 2025, https://developers.cloudflare.com/changelog/2025-04-07-sqlite-in-durable-objects-ga/
10. Durable Objects llms-full.txt - Cloudflare Docs, accessed May 20, 2025, https://developers.cloudflare.com/durable-objects/llms-full.txt
11. Durable Object in-memory state - Cloudflare Docs, accessed May 20, 2025, https://developers.cloudflare.com/durable-objects/examples/durable-object-in-memory-state/
12. Access Durable Objects Storage - Cloudflare Docs, accessed May 20, 2025, https://developers.cloudflare.com/durable-objects/best-practices/access-durable-objects-storage/
13. In-memory state in a Durable Object - Cloudflare Docs, accessed May 20, 2025, https://developers.cloudflare.com/durable-objects/reference/in-memory-state/
14. Cloudflare Durable Objects | TinyBase, accessed May 20, 2025, https://tinybase.org/guides/integrations/cloudflare-durable-objects/
15. Cloudflare Durable Objects - Hono, accessed May 20, 2025, https://hono.dev/examples/cloudflare-durable-objects
16. Getting started · Cloudflare Durable Objects docs, accessed May 20, 2025, https://developers.cloudflare.com/durable-objects/get-started/
17. Durable Object Namespace - Cloudflare Docs, accessed May 20, 2025, https://developers.cloudflare.com/durable-objects/api/namespace/
18. Invoke methods · Cloudflare Durable Objects docs, accessed May 20, 2025, https://developers.cloudflare.com/durable-objects/best-practices/create-durable

-object-stubs-and-send-requests/

19. Cloudflare Durable Objects are Virtual Objects | Lambros Petrou, accessed May 20, 2025, https://www.lambrospetrou.com/articles/durable-objects-are-virtual-objects

20. Thinking in Actors – Challenging your software modelling to be simpler | Hacker News, accessed May 20, 2025, https://news.ycombinator.com/item?id=42291125

21. Durable Object Storage · Cloudflare Durable Objects docs, accessed May 20, 2025, https://developers.cloudflare.com/durable-objects/api/storage-api/

22. Limits · Cloudflare Durable Objects docs, accessed May 20, 2025, https://developers.cloudflare.com/durable-objects/platform/limits/

23. cloudflare-docs/src/content/docs/durable-objects/api/storage-api.mdx at production - GitHub, accessed May 20, 2025, https://github.com/cloudflare/cloudflare-docs/blob/production/src/content/docs/durable-objects/api/storage-api.mdx

24. Use WebSockets - Durable Objects - Cloudflare Docs, accessed May 20, 2025, https://developers.cloudflare.com/durable-objects/best-practices/websockets/

25. Examples · Cloudflare Durable Objects docs, accessed May 20, 2025, https://developers.cloudflare.com/durable-objects/examples/

26. Durable Objects aren't just durable, they're fast: a 10x speedup for Cloudflare Queues, accessed May 20, 2025, https://blog.cloudflare.com/how-we-built-cloudflare-queues/

27. Durable Objects from Cloudflare, accessed May 20, 2025, https://www.cloudflare.com/developer-platform/products/durable-objects/

28. Demos and architectures · Cloudflare Durable Objects docs, accessed May 20, 2025, https://developers.cloudflare.com/durable-objects/demos/

29. Developer Spotlight: Winning a Game Jam with Jamstack and ..., accessed May 20, 2025, https://blog.cloudflare.com/developer-spotlight-guido-zuidhof-full-tilt/

30. Level Up Your Serverless Game with Cloudflare Durable Objects | 9thCO, accessed May 20, 2025, https://www.9thco.com/labs/level-up-serverless-with-durable-objects

31. Ably vs Cloudflare Durable Objects : which should you choose in 2025?, accessed May 20, 2025, https://ably.com/compare/ably-vs-cloudflare-durable-objects

32. Using Durable Objets for Assets - Fast Deployment - Good or bad idea?, accessed May 20, 2025, https://community.cloudflare.com/t/using-durable-objets-for-assets-fast-deployment-good-or-bad-idea/232117

33. Making Super Slurper 5x faster with Workers, Durable Objects, and Queues, accessed May 20, 2025, https://blog.cloudflare.com/making-super-slurper-five-times-faster/

34. Demos and architectures · Cloudflare Pages docs, accessed May 20, 2025, https://developers.cloudflare.com/pages/demos/

35. Guido Zuidhof (Guest Author) - The Cloudflare Blog, accessed May 20, 2025, https://blog.cloudflare.com/author/guido/

36. cloudflare/wildebeest: Wildebeest is an ActivityPub and ... - GitHub, accessed May 20, 2025, https://github.com/cloudflare/wildebeest

37. Tutorials · Cloudflare Durable Objects docs, accessed May 20, 2025, https://developers.cloudflare.com/durable-objects/tutorials/
38. Demos and architectures - Workers AI - Cloudflare Docs, accessed May 20, 2025, https://developers.cloudflare.com/workers-ai/guides/demos-architectures/
39. Tutorial · Cloudflare Durable Objects docs, accessed May 20, 2025, https://developers.cloudflare.com/durable-objects/get-started/tutorial/
40. Error handling - Durable Object - Cloudflare Docs, accessed May 20, 2025, https://developers.cloudflare.com/durable-objects/best-practices/error-handling/
41. Pricing - Workers - Cloudflare Docs, accessed May 20, 2025, https://developers.cloudflare.com/workers/platform/pricing/
42. Pricing · Cloudflare Durable Objects docs, accessed May 20, 2025, https://developers.cloudflare.com/durable-objects/platform/pricing/