

Project Brief: Electron-Based Web Content Capture App (with Playwright & Python)

Project Goals and Scope

The goal of this project is to build a **cross-platform desktop application** that reproduces and extends the functionality of the original Chrome extension ¹ for capturing full web content. Key objectives include:

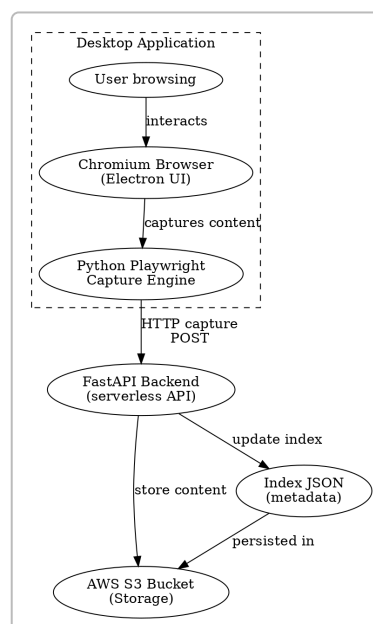
- **Full Content Capture:** Monitor all web page loads in an embedded Chromium browser and capture **everything the page loads**, including the HTML, all JavaScript files, CSS, AJAX/XHR responses, and other resources, in real time. The capture should be **passive and non-intrusive**, meaning it does not alter the target website's behavior or require any modifications on the sites being captured (aside from normal browsing) ². This ensures we archive exactly what the user sees (and what runs behind the scenes) on each page.
- **Preserve and Archive Web Pages:** Store a complete snapshot of each visited page for authorized or target domains (e.g. partner news sites) – including the DOM HTML and loaded scripts – to enable future analysis like fact-checking, provenance verification, or content replay ³ ⁴. Each capture should represent the page as delivered at that time, supporting use cases such as building knowledge graphs or auditing content changes over time.
- **Content Hashing & Deduplication:** Compute a **unique hash** (e.g. SHA-256) of the captured content (combining HTML and script text) to fingerprint each page ⁵. Use this to **avoid storing duplicates** – if an identical page was captured before, the system can detect it via the hash and skip re-uploading, conserving storage and bandwidth ⁶. An index of content hashes will be maintained to prevent redundancy.
- **Client-Side Processing with Python:** Leverage a built-in Python engine to process captures on the client side (the desktop app) as much as possible. This includes hashing, packaging the data payload (e.g. assembling JSON with metadata), and even on-the-fly analysis or content rewriting if needed – all done **locally in Python** for speed and security. By shifting logic to the client, we minimize server load and allow more interactive or real-time processing (similar to the original extension which used Pyodide for client Python ⁷).
- **Serverless Backend for Storage:** Use a minimal **FastAPI** backend (deployable as a serverless function or lightweight service) that simply receives content from the app and stores it in **AWS S3** (or a compatible cloud storage). The backend will remain **stateless** – no databases or persistent state – acting only as a pass-through to storage and an updater for an index file ⁸ ⁹. This keeps the server-side simple, scalable, and reduces security risk (little attack surface beyond the API) ¹⁰.
- **Structured Cloud Storage:** Organize stored pages in S3 with a **timestamp-based key schema**, so each capture is saved under a path by date/time (e.g. `s3://Bucket/YYYY/MM/DD/HH/mm/ss/<page_slug>/page.html`) ¹¹. This provides natural versioning — two captures of the

same URL at different times will reside in different timestamped folders ¹² . Alongside the raw files, maintain a **JSON index** of captures (with fields like URL, timestamp, hash, and file references) to enable quick lookup and to track content evolution over time ¹³ ¹⁴ .

- **Dual Operation Modes:** Support two runtime modes for the app:
- **Production Mode:** The default mode where captured data is **synced to a remote cloud backend** (FastAPI service and S3 storage). This is for live deployment, sending data to AWS (or another cloud) in real time.
- **Development/Offline Mode:** A special mode for developers or offline use, where the app uses a **local FastAPI server** and **LocalStack** (a local AWS emulator) instead of real cloud services. This allows testing and development without internet or cloud costs – all captures are stored to a local S3-mimicking service and processed locally. The functionality remains the same, but everything runs on the developer's machine for quick, cost-free iterations.
- **Cross-Platform Desktop App:** Package the system as a user-friendly desktop application using **Electron**, so it runs on **macOS (primary dev platform)** and is also compatible with **Windows and Linux** from the same codebase ¹⁵ . Electron will embed a Chromium browser instance for the user interface, providing a familiar browsing experience. Users can navigate normally in this app as they would in Chrome, while the app works behind the scenes to capture content. Packaging with Electron also allows native menus, auto-updates, and installers for easy distribution.

In scope for this project is the development of the Electron app, integration of the Python/Playwright capture engine, the serverless FastAPI backend, and the offline mode setup. **Out of scope** are any analysis or UI features beyond the capture itself (e.g. we are not yet building the full analysis pipelines or a complex user interface for viewing archives – those can be future extensions). The immediate focus is a reliable capture and storage pipeline that runs seamlessly as a desktop browser application.

High-Level Architecture and Components



High-level system architecture of the Electron-based capture app. The desktop application contains an embedded Chromium browser (Electron) that the **user interacts with** normally. A **Python engine (Playwright)** runs in parallel to **monitor and capture** everything loaded in the browser. Captured content is sent via HTTP to a **FastAPI backend** (either cloud or local), which stores the data in **S3** and updates an **index** for deduplication and lookup. This design cleanly separates the capture logic, storage backend, and user interface.

The system is composed of several key components working in tandem. Below is a breakdown of each component and their roles:

Electron Desktop Application (Chromium UI)

The desktop client is built with **Electron**, which bundles a Chromium browser engine and Node.js runtime ¹⁵. This allows us to create a native app using web technologies. The Electron app provides the **browser UI** that users will see and use:

- **Embedded Chromium Browser:** The app's main window is essentially a Chromium browser instance. Users can navigate to URLs, click links, and log in to sites just as in a normal browser. We leverage Chromium so that the rendering and behavior of websites is identical to Google Chrome. The user's experience (page layout, media, scripts running) will be the same as a standard browser, ensuring the content we capture is exactly what a real user would see.
- **Navigation Controls and UX:** The app can include basic browser controls (address bar, back/forward buttons, etc.) for usability. On macOS, it will appear as a standard app window. Electron's cross-platform nature ensures the same UI works on Windows and Linux as well ¹⁶. We can package installers for each OS.
- **Integration with Python Backend:** Electron's main process (Node.js) will coordinate with the Python process. For example, on app startup, it can spawn the Python Playwright controller (or establish a connection to it). The Electron side might handle things like a settings menu (e.g. toggle between Production and Offline mode, specify allowed domains to capture, etc.) and pass those settings to the Python layer. Communication can be done through an IPC (Inter-Process Communication) mechanism – such as Electron's `ipcMain`/`ipcRenderer` messages, a local WebSocket, or simply spawning Playwright with the appropriate config.
- **Browser Context Configuration:** Within Electron, we can configure the **browser session** to meet our needs. For instance, we can set a custom user agent if needed, or ensure that the DevTools Protocol is accessible (to allow Playwright to attach). We will also enable proper CORS or permission settings so that the JavaScript within the page (or Playwright script) can send data to our FastAPI backend without being blocked. Electron allows intercepting and modifying web requests (via `session.webRequest` or the DevTools Protocol) if needed, but in our case we mostly rely on Playwright in Python for that logic. The Electron app's role is primarily hosting the UI; it **does not modify page content** except for possibly injecting a content script stub that facilitates communication with the Python side (if needed). The emphasis is that the target websites perceive a normal browser – we don't break any functionality or violate security policies of the sites.

Playwright Python Capture Engine

Instead of running Python inside the browser via Pyodide (as the Chrome extension did), this desktop app runs Python natively on the host machine, using **Playwright** to interface with the browser. **Playwright** is a modern browser automation framework by Microsoft that can drive Chromium (and other browsers) via high-level APIs ¹⁷. We use the **Playwright Python library** to script and monitor the Electron's Chromium instance:

- **Browser Automation & Monitoring:** Playwright connects to the Chromium instance embedded in Electron (likely via the DevTools protocol). It can **control the page** (if needed) and, importantly, **listen to events** like network responses, DOM content loaded, etc. When a page is loaded, Playwright's Python code will capture the **full HTML content** (e.g. using `page.content()` or executing `document.documentElement.innerHTML`) and also gather all **script files and dynamic requests**. Playwright allows hooking into network responses – for every request the page makes, we can intercept the response body. This means even if the page loads additional data via XHR/fetch (AJAX calls), we can grab those JSON payloads or additional HTML fragments in real time ¹⁸. In effect, the Python engine sees everything the browser is loading.
- **Data Packaging and Processing:** Once Playwright collects the raw page HTML and any loaded resources, the Python code will **package the data** for upload. This involves creating a JSON structure containing metadata (page URL, timestamp, hash, etc.) and the content. HTML and text-based resources might be included directly (or Base64-encoded if binary). If a page has many resources (e.g. multiple JS files), the capture logic might bundle them into a single archive (like a zip) or attach them as separate fields. Because this logic is in Python, we have full flexibility in how to structure this payload.
- **Content Hashing & Dedup in Python:** The Python capture engine will compute the **hash of the content** using a library like Python's `hashlib` ⁵. For consistency, we might define the hash input as the concatenation of the page's HTML plus all script texts (or some deterministic combination). This yields a fingerprint we can compare against the index. The Python code can either call a **"check" API** on the backend (e.g. `GET /check?hash=XYZ`) to see if this hash exists ¹⁹, or it might have loaded the index (or recent portion of it) at startup and do a local check. In offline mode, it could directly query the local index file. This deduplication check helps decide if we need to upload the content or skip it. If a duplicate is found, the system might log that "content already archived" and not resend it, or perhaps update a last-seen timestamp in the index.
- **Real-Time Content Rewriting (Optional):** Because Playwright can both **read and inject** into pages, our system can support optional real-time rewriting or annotation of the content. For example, the Python code could inject a small JavaScript snippet into the page after loading, to **highlight certain elements** or to tag resources that were captured. This is optional and would be primarily for debugging or demonstration (ensuring it doesn't disrupt the page's normal functionality). The key idea is we have the power to manipulate the page via Playwright (similar to how a Chrome extension content script could). In normal operation, we likely keep pages unmodified, but this capability could be used for features like showing the user a visual cue that "this page's content has been captured" or masking out sensitive data if needed.
- **Reuse of Python Ecosystem:** By using a native Python environment, we can tap into the full Python ecosystem without the limitations of Pyodide. This means any Python library needed for hashing, data processing, or even performing additional analysis (like parsing text, etc.) can be

installed and used directly. This is a significant improvement in **performance and flexibility**: running Python natively is generally faster and less memory-constrained than via Pyodide in the browser (Pyodide can be 4–8× slower for pure Python code than native execution ²⁰). It also avoids loading a heavy WebAssembly runtime for Python – instead, we use the user’s Python environment bundled with the app. Overall, this makes the capture logic more efficient and easier to develop (debugging Python in a normal runtime, using IDEs, etc., as opposed to debugging in a browser sandbox).

- **“Write Once, Run Anywhere” Logic:** We aim to keep the Python capture code **identical** between the client and server contexts. This echoes the original approach of running the same logic in-browser and on backend ²¹. In practice, this means the Python code that computes hashes and structures data can be reused in the FastAPI backend if needed (for example, to recompute hash on the server for verification, or to perform similar checks). It also means we can write **unit tests** for this capture logic and run them entirely in Python (without needing a browser), by feeding it saved HTML samples. In the original extension, they even ran tests in Pyodide to validate the extension behavior ²²; now we can simply run tests in the normal Python runtime or use Playwright’s headless mode to simulate pages. This alignment of logic ensures consistency and speeds up development: we don’t have one version of code in JS and another in Python – it’s all Python now.

FastAPI Serverless Backend

The backend component remains similar to the original design: a **FastAPI** application that can run as a **serverless function** (e.g. AWS Lambda + API Gateway, or Azure Functions, etc.) or as a lightweight container. Its responsibilities are limited to **receiving content from the desktop app and storing it**:

- **Capture API Endpoint:** The backend exposes an HTTP endpoint (for example, `POST /capture`) that the desktop app will call to upload captured data ²³. The request body is a JSON (or multipart form, if we send files) containing the page content and metadata. Upon receiving a capture, the server will typically **recompute the content hash** (for data integrity and to double-check deduplication) and then proceed to store the files in S3. The backend could also authenticate the request (to ensure only our app or authorized users can send data) – likely via an API key or token that the app includes.
- **Index Management:** If using a JSON index file in S3 to track all captures, the FastAPI function will be responsible for updating this index. For example, after storing a new page, it will append a new entry (URL, timestamp, hash, etc.) to the index file in S3 ^{13 24}. This could involve reading the current index from S3, updating it in memory, and writing it back. Because this is stateless and could be invoked in parallel, we might use S3’s atomic write features or simple locking to avoid race conditions (in low-volume scenarios, simply writing sequentially is fine). In high-volume, a more robust approach (like a DynamoDB index or an append-only log) could be considered, but the scope here assumes moderate capture frequency where a JSON file index suffices ²⁵.
- **Deduplication Check Endpoint (Optional):** We may provide a `GET /check?hash=<hash>` endpoint on the backend that the app can call to quickly ask “Have we seen content with this hash before?” ¹⁹. This would have the backend lookup the hash in the index and return yes/no or the existing record. However, this endpoint is optional; as noted, the app can also just always send captures and let the backend deduplicate by itself (by checking and not re-storing duplicates). An advantage of having an explicit check is saving bandwidth by not sending large

content if we already have it, but the trade-off is an extra round-trip. We will likely implement it for efficiency, especially since checking a hash in a JSON index is a quick operation.

- **Stateless and Scalable:** The backend follows a **stateless, per-request model** – it doesn't keep any session or long-lived state. This aligns with an **"API-first" serverless architecture** as championed by Dinis Cruz and others, where business logic is at the edges (client) and the server mostly acts as a secure storage gateway ⁹. This stateless design means we can scale the capture endpoint easily (multiple instances or lambdas handling requests) and we limit the attack surface. There's no database of sensitive info, and no complex server state to manage ¹⁰. Security is focused on ensuring the upload endpoint is protected (e.g. only our app's authenticated requests are allowed) and that S3 is properly secured. By delegating persistence to S3, we rely on AWS's security for the stored data, which is robust.
- **CORS and Integration:** Since the desktop app is essentially a browser making requests to the backend, we'll configure CORS to allow the Electron app's origin to POST to the API. If the app runs on `http://localhost` in offline mode or a file URL, we might set the backend to be permissive in development. In production, Electron might load `file://` URLs or have an app-specific custom protocol; we'll ensure the backend is accessible to it. FastAPI makes it easy to enable CORS for specific origins.

AWS S3 Storage and Indexing

All captured content will be persisted in **AWS S3**, or a compatible object storage service. S3 serves as a **highly durable, scalable datastore** for the web content:

- **Timestamp-Based Storage Structure:** We use a structured key naming for S3 objects, organizing captures by timestamp. For example, if a page is captured on `2025-03-26` at `16:45:30`, the HTML might be stored at:
`s3://<BucketName>/2025/03/26/16/45/30/<slug>/page.html`
where `<slug>` could be a slugified version of the page title or URL (for readability) or an ID ²⁶. Any additional resources (scripts, images) could either be stored in the same folder (e.g. `.../45/30/<slug>/script1.js`) or bundled into the HTML archive. This timestamped approach ensures that each capture has a unique path (no two captures in the same second collide) and naturally orders the data by time ¹². It effectively treats the S3 bucket as an **append-only ledger** of captured pages – once written, we never modify that capture's file, we only add new ones for new captures, which is good for data integrity and audit trails ²⁷.
- **Index File:** Alongside storing raw content files, we maintain an **index** (likely a JSON file stored in S3, such as `index.json`). This index contains an array of metadata entries for each capture ¹³. Fields in each entry include:
 - `url`: The URL of the page captured.
 - `timestamp`: Capture time (ISO format or epoch) – correlates with the S3 path.
 - `hash`: The content hash fingerprint.
 - `files`: The list of file keys where the content is stored (or a naming convention to derive them).
 - (Optional) `title` or other metadata like content type, etc., if needed for quick reference.

The index allows quick deduplication checks and lookup of content. For instance, to find if we have captured a given URL on a certain date, or to retrieve all versions of a page over time, we can scan or query this index. In a simple implementation, the index is loaded into memory when the app starts (or

when the backend starts) and updated on each capture. For larger scales, we could partition the index by year or use a database, but a single JSON or NDJSON (newline-delimited JSON) file can work for moderate volumes.

- **LocalStack for Offline Storage:** In development/offline mode, S3 calls will be directed to a **LocalStack** instance. LocalStack is a local emulator for AWS services, so it will simulate an S3 endpoint on the developer's machine. The app or backend can be configured (via environment variables or a config file) to use the LocalStack's S3 URL instead of real S3. This provides the same interface (bucket names, object keys) but stores data on the local disk. The benefit is we can test the entire pipeline **without an internet connection or incurring AWS costs**, and with very fast feedback ²⁸. LocalStack's responses are quick (no network latency), accelerating test cycles, and it ensures we don't accidentally pollute real S3 or require cloud credentials for dev. Essentially, offline mode gives us a **sandboxed S3** and a **sandboxed FastAPI** to freely experiment with capturing content.
- **Retention and Volume:** S3 is very scalable, but we will consider how much data we intend to capture. If capturing many pages (e.g. large news sites over time), storage usage can grow. The deduplication via hashing mitigates storing identical content twice. Also, since S3 is pay-per-use, in production we might implement lifecycle policies if needed (for example, archive older captures to Glacier or delete after a certain period if appropriate). These are future concerns; initially, the aim is to retain everything indefinitely to build a comprehensive archive.

Production vs. Offline Mode Details

The application will have a configuration to switch between **Production** and **Development (Offline)** modes, affecting how it connects to backend services:

- **Production Mode:** In this mode, the app uses the **cloud backend** – it will send capture data to the deployed FastAPI service URL (e.g. an API Gateway URL or domain). Likewise, that service will talk to **real AWS S3**. This mode assumes that the user has network connectivity and valid credentials or access for the cloud services. It is the mode for end-users or when running the system in a live environment where data needs to be centralized and persisted in the cloud.
- **Offline Development Mode:** In this mode, everything is self-contained locally. The developer will run a **local FastAPI server** (perhaps the same code as the cloud, just started on localhost), and use **LocalStack** for S3. The app can detect or be started in dev mode (for example, via a command-line flag or an environment variable, or a toggle in a settings file). When in this mode, the Python capture engine might point to `http://localhost:5000/capture` (if FastAPI runs on port 5000) instead of the cloud URL. And the FastAPI in turn will be configured to use the LocalStack endpoint for S3 (which by default might be at `http://localhost:4566` for S3). With this setup, a developer can capture content and verify that everything works with zero cloud usage. It significantly **reduces cost and speeds up testing**, since no real AWS calls are made ²⁸. Moreover, it enables **offline work** – one could develop on an airplane or in a restricted network environment, since all needed services run locally.
- **Parity Between Modes:** We aim for the application to behave the **same way in both modes** from the user's perspective. The only difference is where data is stored. This parity ensures that tests done in offline mode are valid for production. We will likely maintain a single configuration file or set of environment variables that control the endpoints. For example, a config might contain `BACKEND_URL` and `STORAGE_ENDPOINT`. In prod mode, `BACKEND_URL=https://`

`api.mycaptureapp.com` and storage is default AWS; in dev mode, `BACKEND_URL=http://localhost:5000` and storage endpoint overridden to localstack. By keeping these as variables, switching modes is simple and less error-prone.

- **Local AWS Emulation:** Using LocalStack not only avoids costs, but also avoids needing real AWS credentials during development. This is safer (no chance of accidentally pushing test data to production buckets or using wrong credentials). It also allows us to **test edge cases** easily – for instance, we could simulate S3 errors or latency in LocalStack if needed, to see how our system reacts. LocalStack covers many AWS services, so if in future we add, say, DynamoDB or Lambda in our stack, those too can be tested locally. This approach aligns with best practices of using local cloud emulation for rapid, cost-efficient development ²⁹ ³⁰.

Rationale for Using Electron + Chromium + Playwright + Python

This tech stack was chosen to maximize cross-platform support, developer productivity, and to overcome limitations of the previous browser extension approach:

- **Electron (Chromium + Node.js):** Electron enables us to create a desktop app using Web technologies and embed a full Chromium browser ¹⁵. This is ideal because we need a real browser environment to accurately render pages and run their scripts. By using Electron, we **bundle the browser with our app**, ensuring consistency (users don't need to install a specific browser or extension separately) ³¹. Chromium is integrated, so if we ship a certain Electron version, we know exactly what browser engine is running. Electron also provides a way to package the app for **Windows, macOS, and Linux from one codebase** ¹⁶, important for reaching users on all platforms. Another benefit is we can leverage the huge ecosystem of Node.js if needed (for example, using Node APIs for file access, or any of the vast NPM packages for additional features) ³². The choice of Electron also reflects the momentum of web technology: modern web frameworks and dev tools can be used to build our UI, rather than relying on older or platform-specific GUI frameworks ³³.
- **Embedded vs. External Browser:** We specifically embed Chromium (via Electron) rather than requiring an external browser (like launching Chrome with Playwright). This is to provide a **seamless user experience** – the user interacts with our application directly, and we have full control over that browser context. It also simplifies things like making sure the correct extensions or scripts are loaded (since in our own Chromium, we could preload any scripts or set settings as needed). It avoids compatibility issues or user having the “wrong” version of Chrome; everything is self-contained.
- **Playwright for Automation:** Playwright was chosen for the **browser automation and monitoring** piece because it has excellent Python support and can hook into all browser events easily ¹⁷. Compared to writing a Chrome extension (which is constrained by the Chrome Extension APIs), Playwright gives us a **powerful, scriptable control** over the browser. We can intercept network requests, evaluate scripts in page context, and retrieve page state in ways that are more flexible than a content script alone. Playwright also supports **headful mode** (non-headless), meaning we can launch the browser with a UI and still automate it – which is exactly our use case (user-driven browsing with automated observation). By using Playwright's Python capabilities, we keep our core logic in Python rather than JavaScript. This was a conscious choice to unify our implementation language and reuse code from the Pyodide version. Python is also a language well-suited for this kind of data processing and scripting.

- **Python (over JavaScript/TypeScript):** We carry forward the idea of using Python for the core logic (hashing, data processing, etc.) because of its rich ecosystem and the team's expertise. In the original extension, Python ran via Pyodide in the browser; now we can run it natively, which is simpler and faster. Python gives us access to libraries (for example, if we later want to parse HTML or do text analysis, we can use BeautifulSoup, NLTK, etc. with ease). It also means we can share code with any server-side analysis tools in Python. The **performance gain** by using native Python is significant – we avoid the overhead of WebAssembly interpretation (Pyodide) which was reported to be several times slower ²⁰. Additionally, developing in Python outside the browser is easier in terms of debugging and testing (we can run unit tests directly, attach debuggers, profile the code, etc., without the sandbox of a browser).
- **Improved Maintainability and Testing:** Using this stack, we essentially have a **controlled browser environment + an external automation harness**. This is very similar to how automated testing frameworks work, which is a well-understood paradigm. We can write integration tests for our app using Playwright itself (even a script that navigates to a test page and checks that the capture was saved to local storage). We can also easily update the Chromium engine by updating Electron, without worrying about extension manifest changes. The logic running in Python is easier to maintain (one language for logic, rather than split between content script JS and backend Python as before).
- **No Extension Store Limitations:** By moving to a desktop app, we avoid the various **limitations imposed by browser extension platforms**. For example, Chrome extensions (especially under Manifest V3) have restrictions on background scripts, external code execution, and distribution (e.g. need to go through Chrome Web Store or handle manual installs). Our Electron app has none of these restrictions: we can run whatever code we want, use any version of Python packages, and the distribution is under our control (installer or direct download). This also means we can iterate faster without worrying about extension review processes or browser update breaking the extension. We effectively run our own browser instance.
- **Security & Isolation:** Running as a desktop app gives us more control over security. We can ensure that the app only captures content from domains we specify (by coding that rule in our Python logic or Electron config), whereas an extension might risk capturing unintended data if permissions were broad. Additionally, since the entire environment is under our control, we can enforce that no data leaves the app except to our backend. In a standard browser, an extension might be one of many and could conflict or be disabled by the user; in our app, the capture functionality cannot be easily tampered with by third parties. We will, of course, harden the app (enable Electron's security features like context isolation, disable remote module, etc.) to prevent any malicious page from escaping the sandbox or accessing the local system beyond what Chromium normally allows.

In summary, the combination of Electron + Chromium + Playwright + Python was chosen to create a **powerful, standalone tool** that leverages the strengths of each: Electron for cross-platform GUI with a real browser, Playwright for robust control and capture of web content, and Python for advanced processing and easy integration with existing code and tools.

Advantages Over the Original Browser Extension Approach

Transitioning from a Chrome extension to a desktop application brings many technical and practical advantages:

- **Greater Control Over Browser Environment:** The Chrome extension was constrained by the Chrome browser's extension APIs and life-cycle. In the desktop app, we **own the browser instance**. This means we can instrument it at a low level (via DevTools protocol) and are not limited by extension permissions. For instance, a Chrome extension might not easily intercept cross-domain AJAX responses due to CORS or need special permissions; with our own browser + Playwright, we can intercept any request or response programmatically. We can also use debugging techniques (like forcing the browser to enable Network domains, as shown with Electron's debugger API ¹⁸) that wouldn't be possible in a vanilla extension.
- **Elimination of Pyodide Overhead:** The extension ran Python in-browser via Pyodide, incurring performance and complexity costs. In the new approach, Python runs natively, which is faster and uses significantly less memory. We avoid loading a ~ tens-of-megabytes WebAssembly runtime for Python on every browser start. The result is likely a more responsive experience, both for the user (the page isn't bogged down by heavy content script execution) and for the system (processing happens in a separate process). We also avoid the [multi-threading and multi-file limitations of Pyodide](#) (Pyodide runs on the main thread unless explicitly moved to a web worker, which could freeze the UI ³⁴, and had challenges handling multiple files or large packages).
- **Cross-Browser and Future-Proofing:** The extension was tied to Chrome (or Chromium-based browsers). Our Electron app uses Chromium under the hood, but in principle we could switch to a different browser engine or upgrade Chromium by updating Electron. We're not dependent on Chrome-specific extension systems. If in the future we wanted to capture from a different browser (say, an internal corporate browser), Playwright even allows automating Firefox or WebKit with minimal changes. Thus, the architecture is **more flexible**.
- **Distribution and Usage Simplicity:** Some users (especially in enterprise environments) find installing browser extensions difficult due to policies, or they may not use Chrome. Providing a desktop app can be easier in certain cases – users can download an installer and run the app standalone. It also means we can integrate with OS features like auto-launch on startup, dock/taskbar presence, etc., which extensions cannot do. And since Electron apps can auto-update, we have control over pushing updates (without relying on Chrome Web Store updates).
- **UI/UX Enhancements:** As a dedicated app, we can design the UI beyond what a browser + extension can do. We aren't limited to just a popup or devtools panel. For example, we could have a sidebar in the app showing capture status, or a menu to view archived pages, etc. (These are potential future features to leverage the fact this is our own browser UI). The extension approach had a very minimal or no UI (it might have just silently captured or had a simple icon); now we have much more freedom to guide the user or get input (like a list of allowed domains in a preferences screen).
- **Integrated Testing and CI:** We can automate the entire app (using Playwright's ability to launch Electron apps in tests ³⁵) to perform end-to-end tests. This means better reliability. With the extension, testing was trickier – one would need to spin up Chrome with the extension and simulate interactions. Playwright can directly **launch our Electron app** in a test runner and

verify that “if a user visits X page, then Y data is sent to the backend,” etc. This level of testing is easier now, leading to a more robust product.

- **Fewer Security Headaches in Browser Context:** Chrome extensions can be a target for security issues (they have to be careful about cross-site scripting in content scripts, restrictions on eval, etc., and if the extension has broad permissions it's risky). Our app's threat model is a bit different – since it's a browser, it still must safely display arbitrary web content (so we rely on Chromium's security for isolating page content). But we don't have to worry about an external website trying to manipulate our extension via the Chrome extension APIs, because there are none exposed. Also, because the backend is stateless and just a storage proxy, we greatly reduce server-side risk compared to a more complex web application.
- **Customization and Extensibility:** Developing new features is more straightforward in this architecture. Need to capture a new type of resource (say WebSocket messages or video streams)? We can tap into Chromium's devtools events or use Playwright's APIs. Want to change how data is processed? It's Python code we can adjust and test. In contrast, updating the extension's logic might have required dealing with Chrome extension update cycles and limitations on bundling heavy logic (Pyodide was an innovative workaround, but still an extra layer of complexity). Now our stack is basically a normal application – easier to hack on and extend for developers.

In short, moving to an Electron-based app with Playwright gives us **greater power and flexibility** while preserving the core functionality of capturing web content. It addresses many pain points of the extension approach (performance, environment constraints) and opens up new possibilities for the tool's evolution.

Security, Performance, and Extensibility Considerations

Security: The design takes into account several security aspects: - The app will only actively capture content on **whitelisted domains** or user-authorized sites (similar to how the extension only activated on certain sites) ³⁶. This prevents accidentally recording a user's private data from other sites. We will provide a configuration for allowed domains and ensure the Playwright listener only processes those pages. - Data in transit (from the app to the backend) will be sent over HTTPS in production to protect content integrity and privacy. At rest, the S3 bucket can be configured with server-side encryption. Access to the bucket and API will be controlled (e.g. using API keys or signed requests known to the app). - The Electron app will be built with security best practices: enabling context isolation (so that web content cannot access Electron internals or Node directly), disabling remote module, and not loading any remote code except the sites the user browses. Essentially, the embedded Chromium acts like a normal browser – we don't give those web pages any more power than they'd have in Chrome. - The backend being stateless means there's no sensitive info stored there. If an attacker tried to misuse the capture API, they would still need the authorization (which our app holds). We might implement an authentication scheme (like the app signs requests with a secret, or passes a token) to ensure only legitimate clients send data. - One potential concern is that we are capturing and storing potentially large amounts of data, some of which could be malicious (e.g. a page could contain malware scripts). However, since we store content as static files and don't execute them on retrieval, the risk is mainly storage-related. We treat the stored data as untrusted (it's like an archive). Future analysis tools will handle it accordingly (perhaps analyzing in sandboxes). We can also sanitize or exclude certain sensitive data if needed (for example, not capturing cookies or localStorage content that might include user credentials). - If the app has an interface to view or replay captured pages, we would do so safely (e.g.

serve them from a sandboxed viewer, perhaps with scripts disabled to avoid executing old malicious code).

Performance: We are conscious of the performance impact of capturing all resources: - **Async Processing:** The capture engine will run in a separate process (the Python process) so that it does not block the UI. The Chromium browser will load pages normally; Playwright intercepts events asynchronously. This means capturing should not significantly slow down page rendering for the user. There may be minor overhead (e.g. intercepting network responses adds a tiny delay), but we aim to keep it minimal. - We will utilize efficient data handling – for instance, streaming data to S3 instead of building huge blobs in memory if possible. FastAPI and S3 can handle streaming uploads which is useful for very large pages or videos. - The hashing algorithm (SHA-256) is very fast in native code, and even computing it over a few megabytes of HTML/JS is negligible on a modern CPU. So dedup check computations won't be a bottleneck. - We will monitor memory usage. Electron/Chromium is heavy by nature; adding Playwright and Python means extra memory usage. We might allow toggling off capture mode to save resources (if the user wants to browse without capturing for a while). But generally, on a developer machine or modern system, running an Electron app with an extra Python process should be well within acceptable memory/CPU use. - **Network bandwidth** is a consideration: uploading all content of pages (especially if they include large images or videos) could be heavy. In the initial scope, we focus on HTML and scripts (textual content). If needed, we might not capture very large binary resources or we could set size thresholds. Alternatively, we could upload them to S3 but not keep them in memory. We'll design with an eye towards not overloading the user's network – perhaps in a future version letting the user choose to exclude video/media captures. - **Local caching:** We could cache the index or recently seen content list on the client to avoid repeated checks to the server for duplicates (e.g. load the index JSON at startup or have the server return a dedup decision immediately). This would improve performance by avoiding unnecessary uploads or round-trips.

Extensibility: The architecture is modular and **extensible by design**: - The Electron app is essentially a container; we can enhance the UI independently of the capture logic. For example, adding a “*Capture History*” panel that lists recent captures (by reading the index) could be done by reading the index JSON from S3 or a local cache and rendering it in the app. We could also internationalize the UI or add other browser-like features without touching the backend. - The Python capture logic can be extended with plugins for different content types. Currently, HTML and JS are primary, but we could add logic for **PDFs or images**. For instance, if a page loads a PDF (or has an iframe), we might capture that file too. Or if there's an API call returning data, we archive the JSON. We can also integrate processing steps – e.g. after capturing, run the content through a parser or an AI summarizer and store the summary. Because it's Python, integrating with libraries or AI models is straightforward (could be an optional step after storing raw content). - On the backend side, since it's just FastAPI, adding new endpoints or deployment targets is easy. We can deploy the same code to AWS Lambda or run as a Docker container on any cloud. If we later want to support a different storage (say Azure Blob or a Postgres DB), we could implement that behind the same interface. - The use of an **index JSON** is simple but effective. If requirements grow, swapping that out for a more robust database or search index is possible without changing the capture mechanism – it would mostly affect the backend and how the app queries for deduplication. This means our core design (capture everything, store, and index) can scale with different tech under the hood. - **Future integration:** The data we collect could feed into other systems. For example, one goal might be to build a semantic knowledge graph of news content. Since we're storing everything in S3 in a structured way, a separate process or service can ingest those files and perform analysis (e.g. extracting entities, linking stories, etc.). Our capture system lays the foundation for that. By keeping capture/storage separate from analysis, we follow an API-first approach that lets other developers or systems use the archived data easily ³.

In essence, the new architecture not only meets the immediate needs of capturing and storing web content, but is built to be robust, secure, and adaptable for future needs or enhancements. It addresses prior limitations and sets the stage for advanced uses (like integrated analysis pipelines or richer UIs) without significant rework.

Development and Testing Workflow

Developing this project will involve both the Electron app side and the Python backend side. A smooth workflow and robust testing strategy are crucial given the multi-component nature:

- **Project Structure:** We will likely organize the code such that the Electron app (JavaScript/TypeScript code) and the Python capture engine live in the same repository. For example:
 - `electron-app/` - contains the Electron main process and renderer code (the browser window, any UI components).
 - `capture-engine/` - a Python package or script that implements the Playwright logic.
 - `backend/` - the FastAPI application code (which can be deployed or run locally).
 - `config/` - configuration files for modes (could be JSON, or use environment variables).
- We'll include setup instructions for developers to run the app in dev mode.
- **Running in Development:** A developer will be able to run the Electron app in a dev environment (with hot-reload for the UI perhaps) and have it connect to their local Python process. We might create a simple launch script that starts the FastAPI server (for local mode), starts LocalStack (perhaps via Docker if installed, or we instruct devs to start it), then launches the Electron app pointing to local endpoints. This could all be orchestrated with a tool like **npm scripts** or a **Makefile** for convenience.
- **Playwright Integration:** Playwright itself can be used not only in our app but also to test our app. During development, one can run the Python capture engine in a standalone mode for debugging – e.g. have it launch a Chromium instance on its own (outside Electron) to test the capture logic on a sample page. This is useful for verifying that our capture logic works independent of Electron. Playwright has a robust debugging mode where you can see the browser as it runs the script, which helps in writing the interception logic.
- **Unit Testing:** For Python code, we will write unit tests (using `pytest` or `unittest`) for functions like hashing, packaging, etc. These tests can run without any browser – feeding dummy HTML strings and verifying hashes are computed correctly, JSON payloads formed correctly, dedup logic behaves, etc. We can also simulate the index JSON in tests to ensure our duplicate detection works.
- **Integration Testing (Offline mode):** A great advantage of the offline mode is we can do full end-to-end tests **locally**. For example, using a CI pipeline or a local test, we could:
 - Start the FastAPI server locally (on a test port).
 - Start LocalStack (maybe in CI we use the `localstack` Docker image).
 - Run the Electron app in a test mode using Playwright's Electron testing capabilities ³⁵ or via a script. Playwright (or even Spectron, an Electron test tool) can automate the Electron app: open a page, simulate a user clicking a link or entering a URL.
 - The test then waits for the capture to happen (maybe the app emits an event or we check the local S3 for a new file).

- Assert that the file was indeed created in LocalStack's S3 (we can use LocalStack's API or AWS CLI in local mode to list objects in the bucket).

This kind of end-to-end test would validate the whole pipeline in one go. We could have a known small HTML page (maybe served by a simple HTTP server during the test) as the target, so we know what content to expect. After capture, we retrieve the object from S3 and compare it to the source content to ensure fidelity.

- **Playwright Test Scenarios:** We will also use Playwright to test various scenarios like:
 - Navigating to a page that was already captured (should the system skip uploading? We can check that a duplicate wasn't stored again).
 - Capturing a page with multiple resources (ensure all are captured).
 - Behavior when network is down (simulate backend unreachable and see if the app queues the upload or logs an error gracefully).
 - Multi-page navigation (user goes from one page to another quickly; ensure the first page still gets captured).

Playwright's ability to script the browser and interact with our UI is invaluable here. We might write a suite of Playwright tests in Python or TS.

- **LocalStack in CI:** We can integrate LocalStack into our continuous integration environment to test cloud interactions. For instance, spin up LocalStack as a service in GitHub Actions or GitLab CI, run our FastAPI tests against it. The AWS Prescriptive Guidance notes that using LocalStack speeds up tests and avoids costs ²⁸, which aligns with our approach. We will ensure our backend code (FastAPI) has proper abstractions or settings to point to a custom S3 endpoint (LocalStack) when needed.
- **Backend Testing:** The FastAPI backend will be tested separately as well. We can write tests that call the `/capture` endpoint with sample data (bypassing the app) to verify it stores to S3 and updates the index correctly. Using AWS's **Moto** (a library to mock AWS services in Python) or LocalStack in tests can help simulate S3 calls in memory. Since the backend logic is simple (no complex business rules, mostly just storage), testing is straightforward: ensure that given an input, the correct S3 calls happen and the index file ends up with expected content.
- **Developer Environment:** We will provide documentation for setting up a dev environment. Likely prerequisites: Node.js (for Electron), Python 3, and maybe Docker (for LocalStack). The instructions will cover how to run the app in offline mode, how to deploy the backend (for production, e.g. using AWS SAM or Serverless Framework), and how to run tests. We'll also encourage developers to use tools like the Chrome DevTools (accessible via Electron's menus) to inspect pages during development, and the Playwright Inspector for debugging automation scripts.
- **Continuous Integration:** Every commit can trigger a CI pipeline that runs our test suite. This ensures that changes to any part (Electron UI or Python code) don't break the capturing functionality. We might also build the Electron app for all three platforms on CI to catch platform-specific issues (though heavy UI tests could be run on one primary platform like Linux for speed, and just do basic build checks on others).
- **Version Control & Releases:** We will version the project (e.g. v0.1, v0.2, etc.) and use Git for source control. Releases of the Electron app can be packaged with something like Electron Forge or Builder to generate installers. Since the question context suggests an evolving project (the

original was v0.x), we'll maintain a changelog noting the move from extension to desktop app and new features like offline mode.

By following this development and testing workflow, we aim to catch issues early and ensure that the application is reliable. The ability to test offline, and to automate the entire capture process in a controlled environment, will give confidence that the system works as expected before deploying it. Moreover, it sets up a pattern where future contributions (perhaps open-source, as indicated) can be validated with the same rigor. The end result will be a well-tested, **comprehensive web content capture tool** that can be used for archiving and analysis tasks with a high degree of trust in its correctness and performance.

Sources: This project builds upon concepts from the earlier Chrome extension approach ¹ ³ and leverages established tools and patterns in web automation and serverless design. The rationale for an API-centric, serverless backend echoes Dinis Cruz's recommendations for content capture systems ⁹, and the use of LocalStack for offline testing aligns with best practices for cloud development ²⁸. By integrating Electron and Playwright, we combine the proven cross-platform capabilities of Electron ¹⁵ with Playwright's modern automation framework ¹⁷ to achieve a solution that is powerful, flexible, and geared towards future enhancements. The result is an architecture that is both **innovative and pragmatic**, using the right tools for each part of the problem while maintaining clarity of separation between the user-facing app, the automation logic, and the backend storage.

¹ ² ³ ⁴ ⁵ ⁶ ⁷ ⁸ ⁹ ¹⁰ ¹¹ ¹² ¹³ ¹⁴ ¹⁹ ²¹ ²² ²³ ²⁴ ²⁵ ²⁶ ²⁷ ³⁶ [files.diniscruz.ai
https://files.diniscruz.ai/github/pdf/2025/05/18/project_web-content-capture-extension-with-pyodide-and-serverless-backend.pdf](https://files.diniscruz.ai/github/pdf/2025/05/18/project_web-content-capture-extension-with-pyodide-and-serverless-backend.pdf)

¹⁵ ¹⁶ [Build cross-platform desktop apps with JavaScript, HTML, and CSS | Electron
https://electronjs.org/](https://electronjs.org/)

¹⁷ [GitHub - microsoft/playwright](https://github.com/microsoft/playwright): Playwright is a framework for Web Testing and Automation. It allows testing Chromium, Firefox and WebKit with a single API.
<https://github.com/microsoft/playwright>

¹⁸ [JeffProd | Web scraping by watching requests
https://en.jeffprod.com/blog/2021/web-scraping-by-watching-requests/](https://en.jeffprod.com/blog/2021/web-scraping-by-watching-requests/)

²⁰ ³⁴ [Question: How does the speed compare to plain old javascript for regular javascript type stuff? · Issue #776 · pyodide/pyodide · GitHub
https://github.com/iodide-project/pyodide/issues/776](https://github.com/iodide-project/pyodide/issues/776)

²⁸ ²⁹ ³⁰ [Test AWS infrastructure by using LocalStack and Terraform Tests - AWS Prescriptive Guidance
https://docs.aws.amazon.com/prescriptive-guidance/latest/patterns/test-aws-infra-localstack-terraform.html](https://docs.aws.amazon.com/prescriptive-guidance/latest/patterns/test-aws-infra-localstack-terraform.html)

³¹ ³² ³³ [Building a deployable Python-Electron App | by Andy Bulka | Medium
https://medium.com/@abulka/electron-python-4e8c807bfa5e](https://medium.com/@abulka/electron-python-4e8c807bfa5e)

³⁵ [Testing Electron Apps with Playwright - DEV Community
https://dev.to/kubeshop/testing-electron-apps-with-playwright-3f89](https://dev.to/kubeshop/testing-electron-apps-with-playwright-3f89)