**ChatGPT**

# Structure-Aware HTML Diff Algorithms and Implementations

## Challenges in Comparing HTML Documents

- **Preserving Hierarchical Structure:** A robust HTML diff must respect the DOM tree hierarchy instead of treating the HTML as plain text. Traditional line-by-line diffs often break HTML tags or misalign content. Structure-aware algorithms parse the HTML into a DOM tree and compare nodes so that changes are reported at the element or block level (e.g. paragraph, list item) rather than as arbitrary text edits [1]. This prevents invalid HTML in the diff and provides context (for example, recognizing when an entire `<div>` section moved or changed).

- **Ignoring Insignificant Differences (Semantic Equivalence):** Not all textual differences in HTML are meaningful to the rendered content. Effective HTML comparison filters out changes that don't affect the document's semantics or appearance. This includes things like whitespace or line-break variations, different attribute ordering, or changes in tag case [2]. By parsing HTML into a normalized structure and applying rules, diff tools can ignore cosmetic changes (e.g. reordering of `class` attributes or minor formatting) [2]. This ensures the diff highlights actual content changes, not noise.

- **Unified, Human-Readable Output:** Instead of a side-by-side diff, many HTML diffs produce a single **merged view** of the new document annotated with insertions, deletions, or modifications. This "track changes" style output lets users read the latest version with changes marked inline (e.g. added text highlighted or removed text shown with strikethrough). For example, DaisyDiff's output is the newer HTML with special `<ins>` and `<del>` annotations and CSS/JS that highlight added/removed sections [3]. The goal is a self-contained HTML diff that can be opened in a browser to visualize edits in context, much like how Google Docs or Word show revisions within the document.

- **Tracking Moved or Modified Blocks:** A sophisticated diff should recognize when large blocks of content have been moved or substantially rewritten. Pure text diffs usually show a move as a deletion in one place and an insertion in another, losing the connection. Structure-aware algorithms try to maintain **element identity** across versions – e.g. matching the same `<section>` or list item in old and new trees – so that content relocations can be detected. Some diff algorithms (dating back to Heckel's 1978 method) specifically detect block moves [4]. For instance, the wikEd diff library (used for Wikipedia) implements Heckel's algorithm and can flag moved blocks, indicating their original position in the new document's diff output [5]. This way, a paragraph moved from top to bottom is shown as "moved" rather than as unrelated delete/insert chunks.

# Tree-Structured HTML Diff Algorithms

Early research on computing differences in tree-structured data laid the groundwork for HTML-aware diffs. A classic example is the algorithm by **Chawathe et al. (1996)**, *"Change detection in hierarchically structured information"*, which describes how to find changes between two trees (like DOM trees) by identifying corresponding nodes and editing operations [6]. In essence, these algorithms treat an HTML document as an ordered tree and try to find an optimal "edit script" (sequence of insert, delete, update operations) to transform the old tree into the new tree.

A key step is **tree matching:** figuring out which nodes in the original HTML correspond to the same nodes in the modified HTML. Naively matching every node to every other is impractical, so algorithms use heuristics like node identity or similarity. For example, the X-Diff algorithm for XML/HTML uses a *node signature* (based on the node's name and its context path) to constrain matches – e.g. a `<title>` node under a `<header>` won't be matched with a `<title>` under a different section [7] [8]. In general, nodes are only matched if they have the same tag name and play a similar role in the structure (and often similar content). By preserving parent-child relationships in the match, the diff respects HTML hierarchy (a matched subtree in the old and new DOM implies that content is "the same" element before and after, possibly with changes inside).

Once nodes are matched, a **tree diff algorithm** computes what changed. Typical operations include: inserting a new element (and its subtree), deleting an element, moving an element to a new parent, or modifying text content or attributes. Some algorithms do not explicitly have a "move" operation but can infer moves as a delete+insert of a matched subtree. Others, like those based on Heckel's method or similar approaches, will explicitly label a block move [4]. The end result of a tree diff is often an *edit script* or a modified DOM: for instance, X-Diff produces a script of insert/delete/update operations, whereas DaisyDiff directly produces an annotated HTML DOM highlighting the changes.

**Performance:** Finding the minimum edit script for trees is NP-hard in general, so practical implementations use simplifications or cost heuristics [9]. For example, an algorithm might prioritize matching whole moved sections even if not strictly minimal, because that yields a more intuitive diff. Many tools also utilize unique identifiers (like `id` attributes in HTML) to assist matching whenever possible. The outcome is not guaranteed to be the mathematically smallest set of changes, but rather a **meaningful** set of changes that align with human perception of edits [9]. For instance, adding a new column to a table might technically change every row in HTML, but a structure-aware diff would ideally report "added column X" instead of listing many cell-level changes.

After the tree-level diff identifies which elements correspond and which are new or removed, the next step is often an **inline text diff** for the content of changed elements. If a paragraph `<p>` element is matched between versions (same paragraph, but text changed), the diff algorithm will perform a text comparison of the paragraph's text nodes to highlight word-level or character-level changes. Many HTML diff implementations use classic text diff algorithms (like Myers' longest common subsequence or the **diff-match-patch** library) on the text within matching elements. For example, DaisyDiff will drill down into a text node and wrap added words in `<ins>` and deleted words in `<del>` so that changes inside a paragraph are precisely marked [10]. This combination of **tree diff for structure** and **text diff for content** handles cases where content changes span multiple inline elements or vice versa. It ensures that if, say, a sentence was split into two `<span>`s in the new version, the diff can still align the continuous text properly rather than treat it as a completely different structure.

## Semantic Equivalence and Ignoring Trivial Changes

An important feature of "semantic" diff tools is the ability to ignore changes that do not alter the meaning or appearance of the document. HTML files often have irrelevant differences due to editors or formatting – for example, one version might have different indentation, line breaks, or attribute ordering that doesn't affect the rendered page. A naive diff would flag all these changes, obscuring the real content edits.

Structure-aware diff tools mitigate this by normalizing or filtering out such differences. One strategy is to parse the HTML and then serialize or compare it in a normalized form (canonical HTML). Another strategy is to explicitly check potential changes against rendering rules. The **SemanticDiff** tool, for instance, parses HTML into a DOM and then *"applies a set of rules to filter out changes that do not modify the rendering of the document,"* such as ignoring extra whitespace, differences in tag casing, or reordering of CSS classes in a `class` attribute [2] . Only changes that would be visible or meaningful to the user (text changes, content additions/removals, attribute changes that affect display, etc.) are highlighted.

This approach handles issues like inessential whitespace or insignificant attribute shuffling. In practice, many diff libraries offer options to ignore whitespace and line-break differences. The Python `xmldiff` library, for example, allows configuring how to treat whitespace or insignificant text nodes [11] [2] . Similarly, older tools like Ian Bicking's `htmldiff.py` incorporated improvements to skip differences in script tags or whitespace that don't matter [12] .

However, determining semantic equivalence can be tricky. Certain attribute changes (like a change in the `class` or `id` value) might not change the visual output at all, or might trigger subtle CSS differences – the diff tool cannot fully know the impact without CSS context. As noted in a W3C discussion, an inline diff of an attribute change (for example, a `<p class="c1">` changed to `<p class="c2">` ) may confuse users if that class doesn't affect visible output [13] . Many diff interfaces choose not to show such a change inline; or they present it as a side note or require a side-by-side view for non-content changes. The guiding principle is to avoid overwhelming the user with changes that don't affect the reading of the document. By focusing on true content differences and significant structural edits, a semantic HTML diff yields a much cleaner, more **human-readable diff**.

## Presenting Differences in a Unified HTML View

Most HTML diff tools strive to produce an output that is itself an HTML document displaying the changes in context. The unified diff view typically looks like the new version of the document, augmented with visual markers (often colored highlights, inserted text, and struck-out deletions) to denote what changed relative to the old version. This is analogous to reviewing a document with *track changes* enabled.

Technically, the diff output is achieved by injecting special markup into the new document's HTML. A common convention (also used in the W3C HTML diff service and others) is to wrap added content in an `<ins class="diff">…</ins>` tag and removed content in a `<del class="diff">…</del>` tag, often with CSS classes indicating the type of change (insert, delete, modify). For example, if a sentence was removed, the unified diff may still include that sentence in the flow, but wrapped in `<del>` so that it can be rendered with a strikethrough or red highlight; inserted text will appear in an `<ins>` tag with underlining or another color. This way the context is preserved – you read the document normally, with deletions crossed out and additions underlined.

Tools like **DaisyDiff** output the entire new HTML with such annotations. The DaisyDiff library specifically *"annotates a newer HTML document with differences to an older HTML document."* The result contains the full HTML structure of the new version, plus extra markup and CSS hooks to visualize inserted, deleted, and changed elements [3] . Additional JavaScript can enhance the experience (for instance, DaisyDiff's demo page uses JS to navigate between changes and show tooltips with details [14] ). The crucial point is that the diff output remains a valid HTML file that one can open in a browser to see a highlighted, **single-document diff**.

This unified approach is excellent for readability because it mimics how word processors show changes. Microsoft Word's *Compare Documents* feature, for example, generates a single document with revision marks (insertions and deletions in red/green) – essentially the same concept applied to Word's internal XML format. Similarly, Google Docs' *Compare* function or Suggestion mode will highlight additions and deletions in one document view. These systems rely on their internal structured representation to align content and then output the differences as colored highlights in the document. The HTML diff algorithms we discuss operate on the same principle for web documents.

One challenge with unified HTML diffs is ensuring the injected `<ins>` / `<del>` tags (or equivalent markup) don't break the HTML structure. The diff algorithm must insert these tags carefully – usually only at points where element boundaries allow it – or use wrapper elements with special classes if inserting an actual `<del>` would violate HTML nesting. Some older diff tools struggled with this; for instance, one Python HTML diff script naively wrapped a block-level tag in `<del>...</del>` which resulted in invalid HTML [15] . Modern implementations handle this by working with the DOM: they might create placeholder nodes for deleted content, or split nodes as needed to insert diff markers without invalid nesting. The end result should be an HTML file that renders the new document with visual cues for removals and insertions.

It's worth noting that unified diffs are meant **for human consumption**, not as a machine-mergeable patch. They sacrifice exactness for clarity. DaisyDiff's documentation explicitly states that the output will *"not reflect the changes exactly"* in terms of a code diff – some minor content might be omitted or the HTML may be cleaned up in the process – but the goal is to make it easy for a person to see what changed [16] . In summary, the unified diff format presents a *merged view* of the two versions, highlighting changes at the block and inline level, and is the format used by most HTML-aware comparison tools and word processors.

## Detecting Moved Content and Maintaining Identity

Detecting moved blocks is an advanced aspect of diff algorithms. In the context of HTML, a "move" means a chunk of the DOM (one or more nodes and their children) that existed in the original document now appears in a different place in the new document. A naive diff would treat this as a deletion (from the old location) plus an insertion (at the new location). However, users often want to know that it's the *same* content just relocated, especially for large sections.

Structured diff algorithms can detect moves by leveraging the matching phase. If the algorithm can match a node in the old tree with an identical node in the new tree *despite* it having a different parent or position, then it has essentially found a move. For example, if an entire `<div id="sidebar">` block moved from the left column to the right column in a page, a good matching algorithm will still pair the old and new `<div id="sidebar">` nodes (since their id and content are the same), rather than treating them as

unrelated. The resulting edit script might then represent this as a "move operation" or as a deletion+insertion but with an annotation that they are the same node.

Historically, one of the first algorithms to handle move detection in diffs was Paul Heckel's technique (1978) [17] [4] . Heckel's approach was linear-time and could identify block moves when there weren't too many duplicate lines. Modern variations of this are used in code diffs and have been adapted to HTML/text. The **wikEd diff** JavaScript library, for instance, implements Heckel's algorithm and extends it to output a rich diff with moved blocks highlighted [4] . In a visual diff, a moved section might be indicated by a special color or annotation (e.g. "moved from earlier in the document") rather than just a simple insert. WikEd diff even shows the original position of the moved text in the output [18] , which helps the user understand relocations.

For HTML, maintaining identity might also be aided by attributes. If elements have stable IDs, a diff tool can use those as anchors (i.e. if `<section id="intro">` exists in both versions, they're likely the same section even if order changed). Some XML diff tools allow specifying a unique key for matching nodes (for example, matching table rows by a key column). The Python `xmldiff` library supports user-defined *unique attributes* to improve matching [11] , which can be very useful for HTML (e.g., matching `<tr>` by a primary key contained within, etc.). In absence of explicit keys, algorithms may fall back to content similarity – for instance, matching a paragraph that moved by its text content.

In practice, not all HTML diff implementations try hard to detect moves, because it can be complex and sometimes unreliable if content is edited during the move. Some tools simply report moves as delete/insert pairs. But higher-end document comparison software (such as Microsoft Word's compare feature) does detect moves: Word will mark moved text with a special double-strikethrough (at the original location) and double-underline (at the new location) to show it as a move. This is essentially the application of a structure-aware diff on the document's XML with additional logic to label moves when the deleted text exactly matches inserted text elsewhere. We see similar capabilities in code-oriented semantic diffs (e.g. SemanticDiff for code can detect refactoring moves). For an HTML diff tool intended for documents or web content, move detection is a desirable feature, and it relies on the underlying tree matching algorithm being clever about pairing up identical or highly similar subtrees across versions.

## Notable Algorithms and Tools for HTML Diffs

Over the years, many libraries and algorithms have been developed to perform structure-aware diffs on HTML or XML:

- **XML/HTML Tree Diff Algorithms:** Research-driven algorithms like Chawathe's tree diff (1996) and subsequent XML diff algorithms (XyDiff, X-Diff, etc.) laid out methods to compute changes between tree-structured data. These algorithms aim to find an optimal or near-optimal edit script and often serve as the engine inside diff tools. For example, the open-source Python `xmldiff` library implements a tree comparison approach to generate human-readable diffs for XML/HTML content. It emphasizes that *"change detection in hierarchical data is very different from flat data"* and provides tools to create diffs that users can read, rather than low-level line diffs [11] . Such libraries typically output changes as a sequence of operations (insert node X, delete node Y, etc.) or allow formatting the diff in various ways (plain text, XML patch, or even an HTML visualization).

- **Early HTML Diff Utilities:** A number of earlier tools were created to specifically compare HTML while respecting tags. For instance, Bert Bos's **HtmlDiff 0.4** (written in C) was known to be fast, though it struggled with very large or radically different documents [19] . On the other end, a Python script by Ian Bicking (also called **HtmlDiff**) handled even radical HTML changes but was *"quite slow for large files"* [19] . Aaron Swartz wrote another HTML diff in Python (GPL-licensed), which was noted to be **unusably slow** on large inputs [20] . On the Perl side, there's **HTML::Diff** and Norm Walsh's **diffmk** tool. Diffmk was used in early W3C HTML diff services – it would take two HTML files and produce a combined file with change markers. The W3C's online HTML Diff service (circa 2007) was based on a Perl script and an alternate version of it was published with source code [21] . These early tools demonstrated the concept, though often had performance or output-format quirks. A W3C wiki comparison noted, for example, that Bicking's Python htmldiff would catch a changed attribute but output it in a not-fully-usable way (inserting raw `<del>` tags in the markup) [15] .

- **LXML and Other Modern Libraries:** With the advent of powerful HTML parsers, some libraries integrated diff capabilities. Python's **lxml** library (for XML/HTML processing) includes `lxml.html.diff` , which *"seems to work pretty well"* according to users [22] . This function can produce an HTML diff by leveraging lxml's parsing of the DOM. Similarly, Java's Apache XML libraries and others have XML diff tools that can be applied to HTML (provided the HTML is well-formed or can be tidied to XHTML). These tend to use tree-diff algorithms under the hood and output either an edit script or an annotated document.

- **DaisyDiff (Java):** One of the most well-known structure-aware HTML diff tools is **DaisyDiff**, an open-source Java library originally released in 2007 by Google. DaisyDiff parses the HTML into a DOM and then compares nodes hierarchically. It *"highlights added and removed words and annotates changes to the styling"* within the document [23] . Notably, DaisyDiff doesn't simply mark an entire changed element as added or removed; it can pinpoint the specific text that changed inside a block, making the diff output very granular and reader-friendly [10] . The output is an HTML file with inserted `<ins>` / `<del>` elements and CSS classes (the library provides a CSS file for coloring and can use a bit of JS for interactive features like tooltips) [14] [3] . DaisyDiff became popular and was even integrated into MediaWiki for a time to provide "visual diffs" on wiki pages. (It was later removed due to some bugs, but it spawned derivatives in other languages.) There was also a PHP port of DaisyDiff (sometimes called VisualDiff, used in early versions of MediaWiki's Visual Editor) [24] . DaisyDiff's algorithm tries to preserve the HTML structure rigorously (ignoring insignificant differences) so the resulting diff HTML is clean and easy to navigate [1] . Users have praised its clear visualization of changes, though some noted it had difficulty with certain cases like identical images being flagged as changed [25] .

- **HtmlDiff in Other Languages:** The idea of DaisyDiff/HtmlDiff has been reimplemented in various languages. For example, a C# port of a Ruby HTML diff library was created (sometimes referred to as **HtmlDiff.NET**), and subsequently a JavaScript port **htmldiff.js** was derived from that [26] [27] . These ports usually follow the same approach: take two HTML snippets, break them into DOM or token streams, find matches, and produce a merged HTML with `<ins>/<del>` tags. The JavaScript implementations (like **htmldiff.js** or similar NPM packages) enable in-browser HTML comparison. They are useful if you want to do a diff on the client side (for example, in a rich text editor on the web) and display the result without a server round-trip. Many of these focus on inline text differences and assume the overall structure is similar; they might not handle large reordering as well as a true tree-diff algorithm, but they are straightforward to use for small content changes.

- **Visual Document Comparison Tools:** Beyond open-source projects, there are commercial or specialized tools aimed at high-quality HTML or document diffs. **HTML Match** (a Windows GUI tool) and **Draftable** (online/API document comparison) are examples that compare HTML documents while preserving formatting. According to one analysis, HTML Match was able to preserve all styling and provided the highest quality visual comparison (essentially rendering the two documents and highlighting changes) [28] . These tools often combine DOM diffing with rendering knowledge to ensure that, say, changes in font or styling are caught. They aren't open-source, but they demonstrate the state of the art: high fidelity comparison that looks close to a manual review.

- **Difftastic (Rust):** In the realm of code and data diffs, **Difftastic** is a new-generation diff tool that can handle structured formats (including HTML) by parsing them. It uses Tree-sitter grammars to parse files and then compares the syntax trees instead of raw text. The result is a diff that *"understands syntax"* and *"produces accurate diffs that are easier for humans to read."* In the case of HTML, Difftastic will parse the HTML into a tree and can therefore detect changes such as added wrappers or moved attributes more intelligently than a line diff [29] . It also explicitly ignores pure formatting changes — for instance, if an HTML tag was just split into multiple lines or spaced differently, Difftastic's output would show no change [30] . While primarily a CLI tool (often used with source control diffs), it illustrates the benefit of syntax-aware diffing for HTML. It does not output an HTML file with highlights (it outputs a diff in its own textual format or side-by-side view), but it's useful for developers reviewing HTML code changes.

- **SemanticDiff and Similar Projects:** SemanticDiff (as mentioned earlier) is another modern initiative focusing on meaningful diffs. It provides a VS Code extension and online service for comparing not only code but also HTML, JSON, etc., ignoring irrelevant differences [2] . In a Reddit post, the author of SemanticDiff mentioned plans to hide style-only changes and even detect moved code blocks for programming languages. This kind of tool often builds on the lessons from both document diffs and AST-based code diffs – meaning it could be using a combination of tree parsing, heuristics for what counts as a real change, and algorithms to detect reordering or moves. The online HTML diff from SemanticDiff shows a side-by-side aligned view, which is slightly different from the unified output approach, but under the hood it aligns the DOM structures and then highlights differences in tags, attributes, and text while collapsing identical sections [31] [2] .

In summary, if you are implementing an HTML comparison in a JavaScript front-end and Python back-end (as many modern web apps might), you have a variety of approaches and libraries to draw from. On the Python side, you could leverage libraries like `xmldiff` for heavy lifting or use an existing HTML diff script (there are Python ports of HtmlDiff). On the JavaScript side, you could use a library like htmldiff.js or even integrate a WebAssembly version of a diff for performance, depending on where you want the computation. The key is to parse the HTML into a structure on at least one side, normalize/ignore trivial changes, then compute differences at the node and text level. By using the techniques and algorithms above – tree-based matching (to preserve structure and detect moves), and inline diff for text – you can present the user with a clear, unified view of changes that feels as if the document had "track changes" on. The combination of academic algorithms (tree edit distance, node signature matching, etc.) with practical heuristics (whitespace ignore, key attributes) is what makes a robust HTML diff.

Finally, observing how **Google Docs** or **Microsoft Word** handle document version comparisons confirms these principles: they operate on a rich internal representation (not flat text) and present differences as marked-up documents with insertions, deletions, and styling changes noted. While their exact algorithms

aren't public, the comparison outputs are essentially what we aim for – a human-readable, semantically meaningful diff of rich text. By using the tools and methods described above, one can achieve similar results for HTML documents, ensuring that structural changes and content edits are conveyed clearly while trivial alterations are filtered out.

**Sources:**

1. Cacciari Miraldo & Swierstra (2017). *"Structure-aware version control: A generic approach using Agda."* – Highlights the limitations of line-based diffs and explores tree-structured diff algorithms [32] [9] .

2. Bos, Bicking, Swartz, Walsh, *et al.* – W3C HTML Diff tool comparisons (circa 2007–2010). Provides insights into early HTML diff tools (C, Python, Perl) and their trade-offs [19] [33] .

3. **SemanticDiff (HTML Online Diff)** – Documentation on structure-based HTML comparison, ignoring non-rendering changes [2] .

4. Wang et al. (2003). *"X-Diff: An effective change detection algorithm for XML."* – Describes using node signatures and tree matching to compare XML/HTML trees [7] [8] .

5. Aucke Bos (2018). *"Visualizing differences between HTML documents"* (Bachelor's thesis). – Introduces the HDiff Java tool based on Chawathe's algorithm and compares DaisyDiff and others. Notably mentions DaisyDiff's approach to inline text differences [10] [23] .

6. **Gentics DaisyDiff Documentation** – Confirms DaisyDiff's tree-based approach and that the output is an annotated HTML for human consumption [1] [3] .

7. Stack Overflow – discussions on diffing HTML and moved blocks. e.g., explanation of Heckel's algorithm and wikEd diff for moved text [4] [18] .

8. Wilfred Hughes – **Difftastic Manual** (2022). Explains syntax-aware diffing (using Tree-sitter for HTML) and ignoring purely formatting changes [29] [30] .

9. W3C HTML Diff service example – Demonstrates output formats and challenges (e.g., attribute change handling) [15] [13] .

---

[1] [3] [16] Gentics Content.Node - REST API Library: DiffResource

https://www.gentics.com/Content.Node/guides/restapi/resource_DiffResource.html

[2] [31] SemanticDiff - Online HTML Diff

https://semanticdiff.com/online-diff/html/

[4] [5] [17] [18] Is there a diff-like algorithm that handles moving block of lines? - Stack Overflow

https://stackoverflow.com/questions/10066129/is-there-a-diff-like-algorithm-that-handles-moving-block-of-lines

[6] [10] [14] [23] [25] cs.ru.nl

https://www.cs.ru.nl/bachelors-theses/2018/Aucke_Bos___4591496___Visulizing_differences_between_HTML_documents.pdf

[7] [8] Microsoft Word - paper401.doc

https://pages.cs.wisc.edu/~yuanwang/papers/xdiff.pdf

[9] [32] webspace.science.uu.nl

https://webspace.science.uu.nl/~swier004/publications/2017-uu-cs.pdf

[11] xmldiff — xmldiff documentation

https://xmldiff.readthedocs.io/

[12] [13] [15] [19] [20] [21] [22] [33] HtmlDiff - W3C Wiki

https://www.w3.org/wiki/HtmlDiff

[24] [26] HTMLDiff Software Projects and Repos (Open Source)

https://onassar.github.io/blogs/web/2012/11/21/htmldiff-software-discoveries/

[27] htmldiff-js CDN by jsDelivr - A CDN for npm and GitHub

https://www.jsdelivr.com/package/npm/htmldiff-js

[28] Baltic Journal of Modern Computing, 2013, Vol.1, No.1

https://www.bjmc.lu.lv/fileadmin/user_upload/lu_portal/projekti/bjmc/Contents/3_3_4_Ostreika1.pdf

[29] [30] Difftastic, a structural diff

https://difftastic.wilfred.me.uk/