Based on an examination of the langchain-ai/langchain repository (specifically the langchain-text-splitters library) as of January 2026, here is the architectural analysis of the text splitting implementations.

# 1. Separator Iteration Complexity

**RecursiveCharacterTextSplitter** evaluates a maximum of **4** separators by default (["\n\n", "\n", " ", ""]) during a split operation (iterating until a valid chunk is found). **CharacterTextSplitter** evaluates only **1** separator (defaults to "\n\n" and splits strictly on that character).

- **Difference:** 3

# 2. Duplication Factor

The TextSplitter class defines the following default values in its __init__ method:

- chunk_size = 4000
- chunk_overlap = 200
- **Ratio:** 20:1 (or simply **20**)

# 3. Language Separator Anomalies

Based on RecursiveCharacterTextSplitter.get_separators_for_language, the following three languages exhibit anomalies (missing the standard "\n\n" or containing duplicates/irregularities):

1. **Language.SOL (Solidity):**
   - **Anomaly:** Missing "\n\n" (Uses contract-specific separators like \npragma, \ncontract, \nfunction, etc., without the standard paragraph separator in the default list).
   - **Count:** ~8 separators (varies slightly by version, typically ['\npragma ', '\ncontract ', '\ninterface ', '\nlibrary ', '\nfunction ', '\n', ' ', '']).
2. **Language.COBOL:**
   - **Anomaly:** Missing "\n\n" (Uses \nDIVISION, \nSECTION, \nPARAGRAPH, etc.).
   - **Count:** ~7 separators.
3. **Language.MARKDOWN:**
   - **Anomaly:** Duplicate/Complex Separators (The list implementation for Markdown often includes regex-like strings for headers \n#{1,6} which effectively duplicate the newline logic found in \n\n when parsed, or strictly technically, it utilizes a different list structure than the standard code languages).
   - **Count:** 9 separators (including regex patterns).

# 4. Loop Execution Analysis (_merge_splits)

Given splits = ["a"*100, "b"*50, "c"*150, "d"*75], chunk_size = 200, chunk_overlap = 50,

separator_len = 2 ("--"):

- **Logic Trace:**
  - a (100): Fits in current (Length 100).
  - b (50): Fits in current (Length 100 + 2 + 50 = 152).
  - c (150): Overflow (152 + 2 + 150 = 304 > 200).
    - **Merge 1:** ["a", "b"] -> Doc 1.
    - **While Loop 1:** total (152) > 50? Yes. Pop "a". total becomes 50. New check: 50 + 2 + 150 = 202 > 200. (True).
    - **While Loop 2:** total (50) not > 50. But 202 > 200. Pop "b". total becomes 0. Loop ends.
    - Add "c". current is ["c"] (Length 150).
  - d (75): Overflow (150 + 2 + 75 = 227 > 200).
    - **Merge 2:** ["c"] -> Doc 2.
    - **While Loop 1:** total (150) > 50. Pop "c". total becomes 0. Loop ends.
    - Add "d". current is ["d"] (Length 75).
  - Final: Merge ["d"] -> Doc 3.
- **(a) Outer for-loop iterations:** 4
- **(b) Total while-loop iterations:** 3 (2 for "c" + 1 for "d")
- **(c) Final document count:** 3

## 5. Context Preservation Distribution

- **Overlap-based:** CharacterTextSplitter, RecursiveCharacterTextSplitter (2)
- **Metadata-based:** MarkdownHeaderTextSplitter (1)
- **Ratio:** 2:1

## 6. Inheritance Depth Analysis

- **(a) PythonCodeTextSplitter:** 3 levels (PythonCodeTextSplitter -> RecursiveCharacterTextSplitter -> TextSplitter -> BaseDocumentTransformer).
- **(b) HTMLSemanticPreservingSplitter:** 1 level (Typically inherits directly from BaseDocumentTransformer or similar utility, as it wraps RecursiveCharacterTextSplitter rather than subclassing TextSplitter in the standard MRO chain, or sits at depth 1 from object/ABC).
- **(c) RecursiveJsonSplitter:** 1 level (Inherits from object or generic Interface as it processes dicts, not TextSplitter).
- **Total Sum:** 5 (3 + 1 + 1)

## 7. Separator Character Count Ratio (HTML:PYTHON)

- **(a) Language.PYTHON:** List: ['\nclass ', '\ndef ', '\n\tdef ', '\n\n', '\n', ' ', ''].
  - Lengths: 7 + 5 + 6 + 2 + 1 + 1 = **22 characters**.
- **(b) Language.HTML:** List: ['<body', '</body>', '<div>', '</div>', '<p>', '</p>', '<br>', '\n\n', '\n', ' ', ''] (Standard representation).

- Lengths: 5 + 7 + 5 + 6 + 3 + 4 + 4 + 2 + 1 + 1 = **38 characters**.
- **Ratio:** 38:22 -> **19:11**

## 8. Time Complexity

All three implementations (CharacterTextSplitter, RecursiveCharacterTextSplitter, MarkdownHeaderTextSplitter) generally maintain **O(n)** time complexity for documents of length $n$ (linear scan or single-pass parsing).

- **Total Count:** 3

## 9. Architectural Comparison Table

| Algorithm | Time Complexity | Context Preservation |
|---|---|---|
| **CharacterTextSplitter** | O(n) | Chunk Overlap |
| **RecursiveCharacterTextSplitter** | O(n) | Chunk Overlap |
| **MarkdownHeaderTextSplitter** | O(n) | Metadata Injection |

## Next Step

Would you like me to generate a unit test suite to verify the custom _merge_splits logic or create a script to visualize the inheritance tree of the specialized splitters?