When exporting crawled documents to a file system, by default, Content Analytics exports content as .dat files *and* preserves the binary format of the file. For example, you export a .doc source file in the .dat format and rename the .dat format to the .doc format. In this case, the resulting .doc file is the same as the original .doc file. However, when you export to a .csv file, the binary content of the document is not exported.

When you configure the export of crawled content for Content Collector integration, Content Analytics preserves the extension from the source (if available) and exports content with the original extension. For example, if the source document has the document sample.doc file, the exported content also has the sample.doc file name. For situations where Content Collector integration is enabled, but the source document does not provide a file extension, the document is exported as a .dat file.

### Binary content file name

When you configure Content Analytics to export content to a file system, the name of the XML file begins with 0000.dat and increments sequentially. For example, if the source file is the sample.doc file, the metadata file name is 0000.dat or 0000.doc if Content Collector integration is enabled.

### Binary content in the relational database

Binary content is stored in the relational database as a binary large object (BLOB).

## 10.3.4  Common Analysis Structure format

CAS is a data structure for representing information that is gathered during the analysis of document such as annotations, tokens, and facets. When you export the CAS format to a file system, the data is exported in the XMI format as .xmi files. The CAS format conforms to UIMA standards and is subject to change with future releases of UIMA.

When Content Collector integration is enabled, the format of the CAS file does not change.

### Common Analysis Structure file name

The name of the XMI file begins with 0000.xmi and increments sequentially until the name reaches 9999.xmi. When more than 10,000 documents are exported, a new folder is created with file names beginning with 0000.xmi.

### Common Analysis Structure in relational database

CAS export to relational database is not supported.