

The global Font Style Gallery contains descriptions of all the font styles (comprised of font types, font styles, sizes and size ranges) in the original document and enables the restoration of its original look and feel as needed. The Style Gallery is common to all Entities and is represented as an external unit in the XML of each print issue remaining unique to the edition, catering to the dynamic nature of print design.

2.6 Image Mapping and Indexing

Effective access and search of historical documents generated from paper and microfilm is an important part of the PrXML effort. Scanned images of such materials are usually degraded and is difficult to process with regular OCR (Optical Character Recognition) technologies.

The readability factor is greatly increased allowing users to access and read the images directly instead reading poorly OCR generated text. The task of comprehending the degraded text is then transferred to the human eye and brain—which is arguably the best OCR.

However what about the searchable factor? How can we search for such images? The unlinking of the 'searchability' factor from the 'readability' is provided by a patented technique called Bitmap Indexing™. This technique is based on the mapping of each meaningful group of pixels (containing a page element like a title, body text, word-patterns or picture) on the page image. Indexing such valuable image mapping data, enables direct search and access to, and sophisticated manipulation of image "clips" instead of cumbersome access to the page image as a whole. Bitmap Indexing results in meaningful end user features, for example, search hits can be highlighted in an article image and search result pages can display scaled images of article titles, not imperfect OCR text.

Olive's PrXML also defines special word patterns tags called APFS (patent Pending) tags (Adaptive Probability Fuzzy Search). These word patterns define the quality parameters and the probability for suspected OCR errors in each word. Smart indexing of these tags enables a search engine to effectively search within degraded images, while compensating on possible OCR mistakes. Without APFS, the regular fuzzy search feature that normally exists in many search engines, will provide a mass of irrelevant search results because the fuzzy search is performed on the entire document text instead on just on the suspect misspelled words.