# Using textual context for improving OCR performance in biomedical literature retrieval
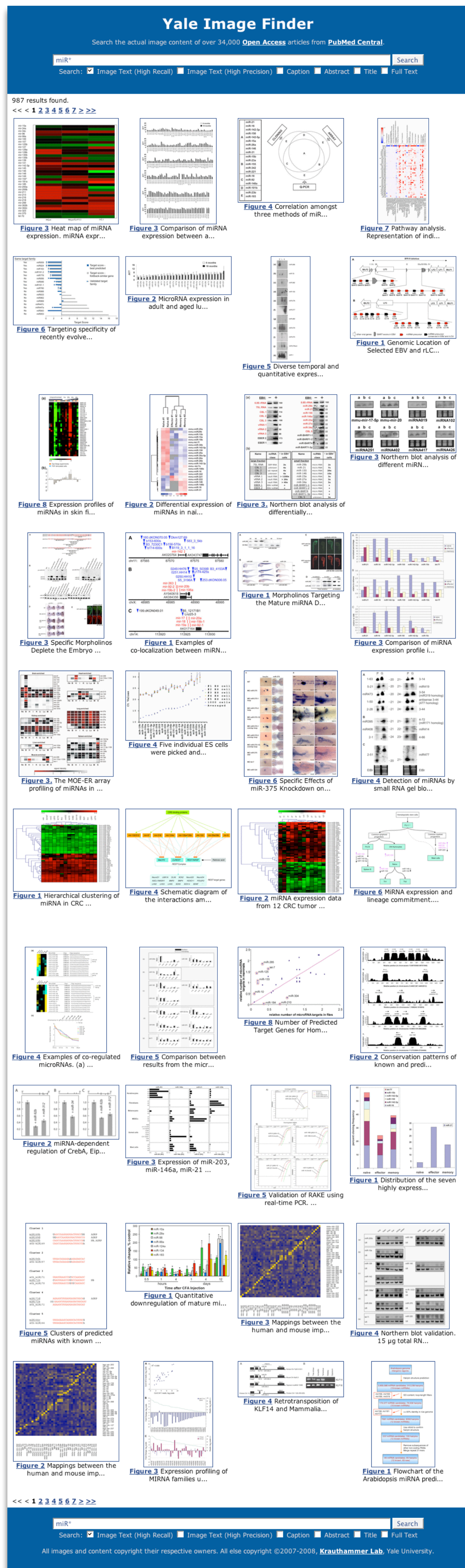
Songhua Xu[1], James McCusker[2], Martin Schultz[1], and Michael Krauthammer[2]

[1]Department of Computer Science, [2]Department of Pathology and Yale Center for Medical Informatics, Yale University, New Haven, CT
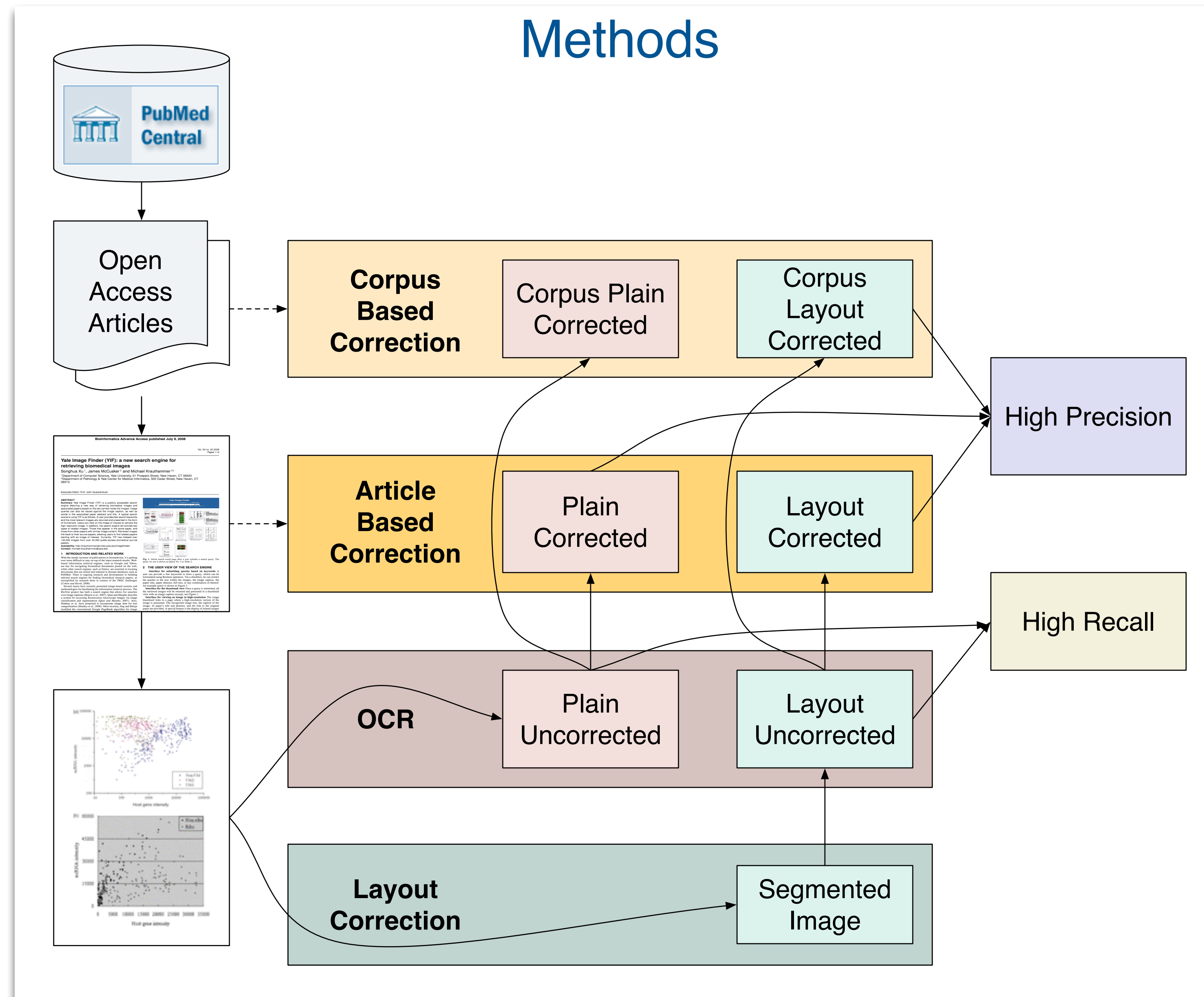
http://krauthammerlab.med.yale.edu/imagefinder

## Abstract

Today's information retrieval (IR) techniques are mostly text-based, which fail in situations when textual information is not easily accessible, such as in biomedical images and figures. We propose to augment IR with optical character recognition (OCR) capabilities, and describe a context-based method for boosting OCR performance.
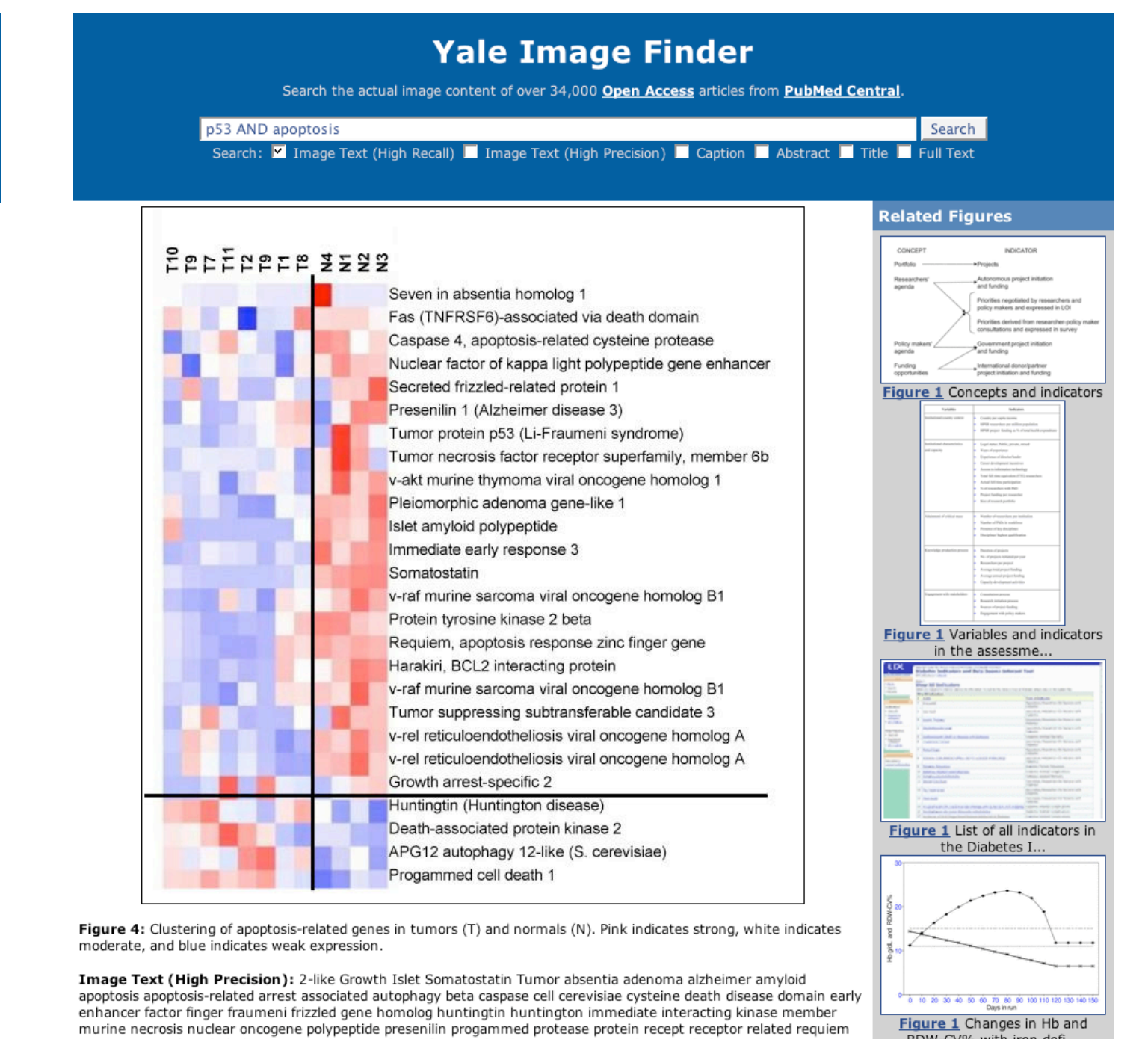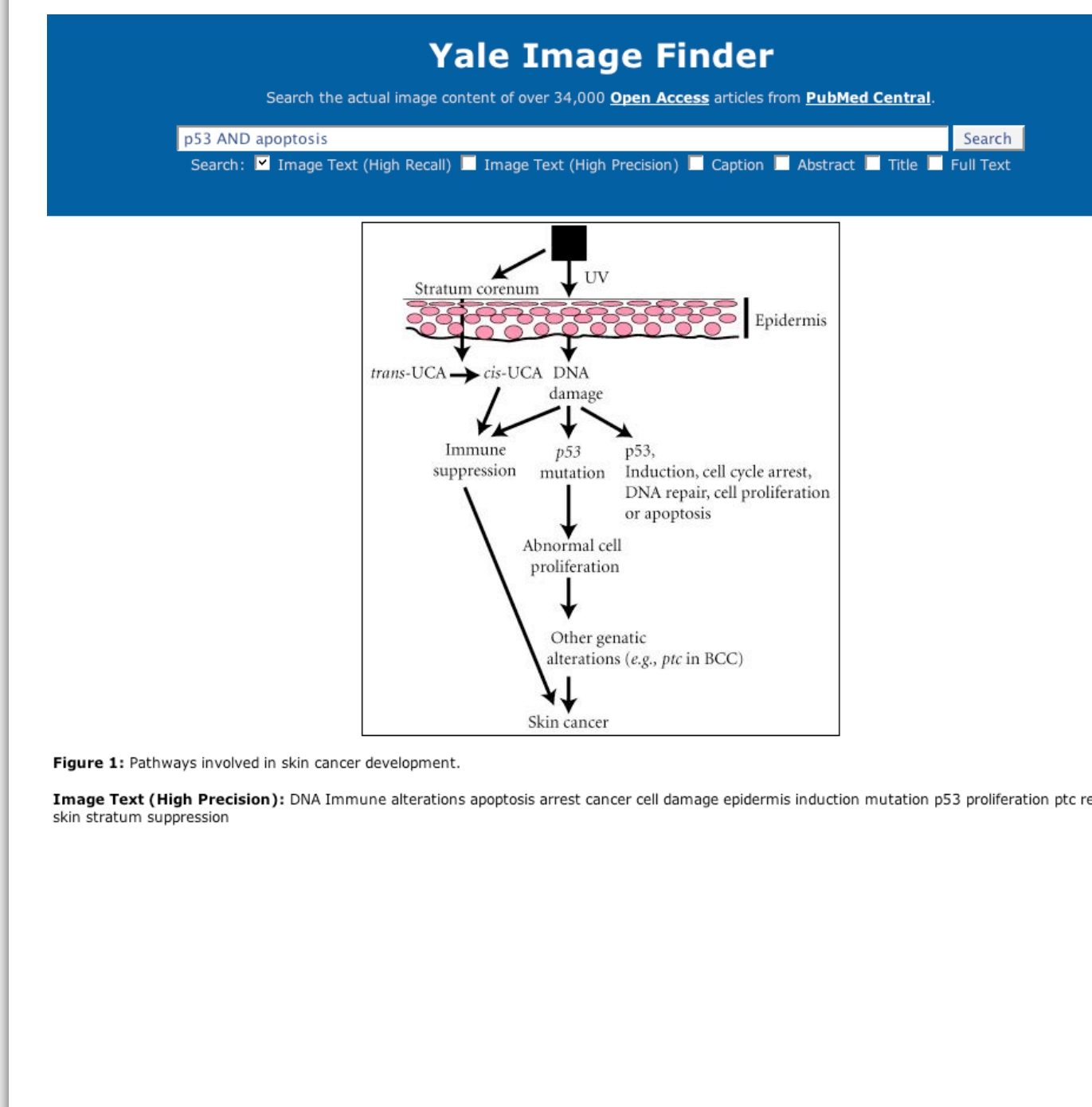
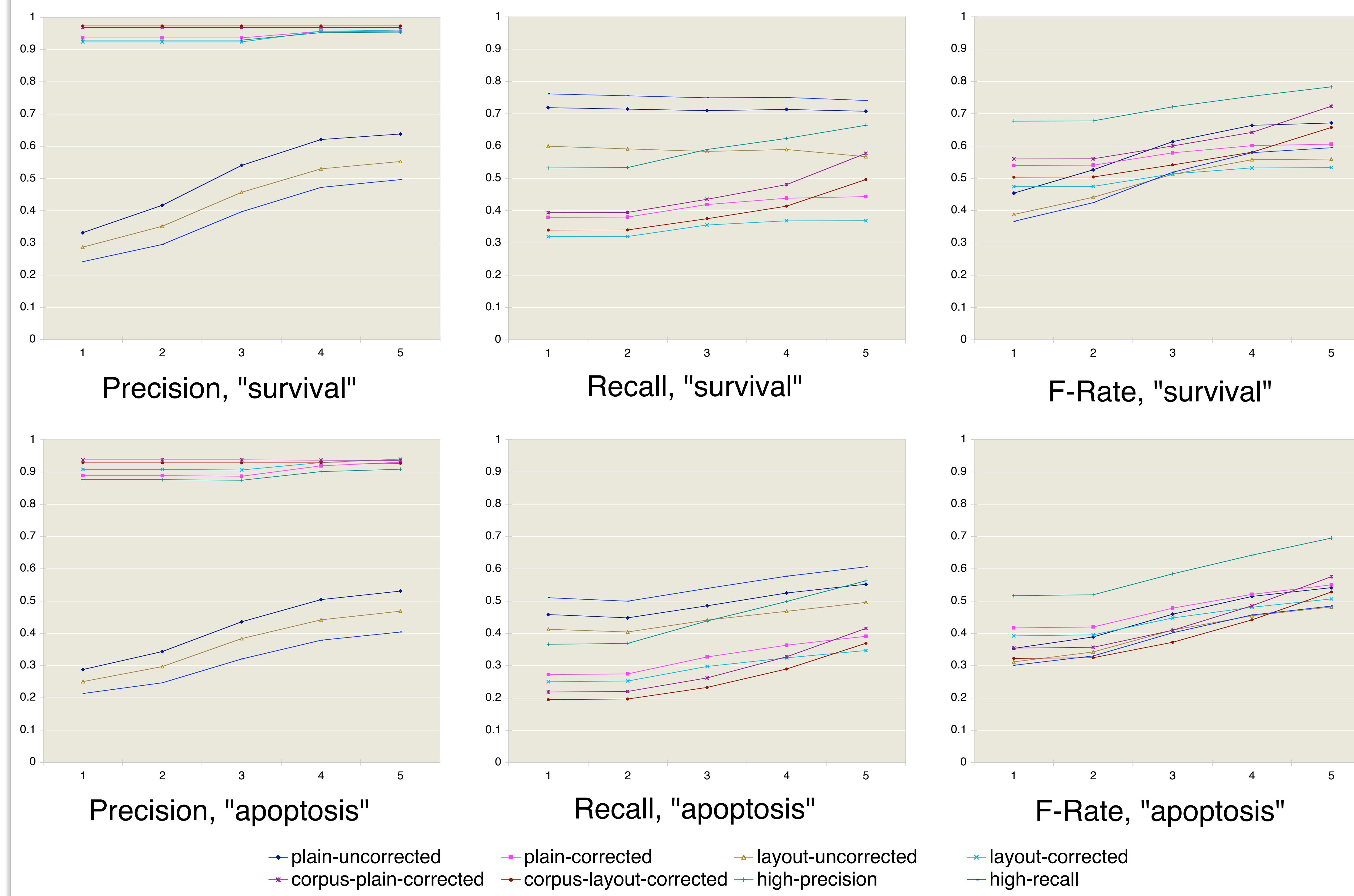Yale Image Finder search results for "miR*" in high-recall mode.

## Methods



## Advantages of OCR-based image text extraction

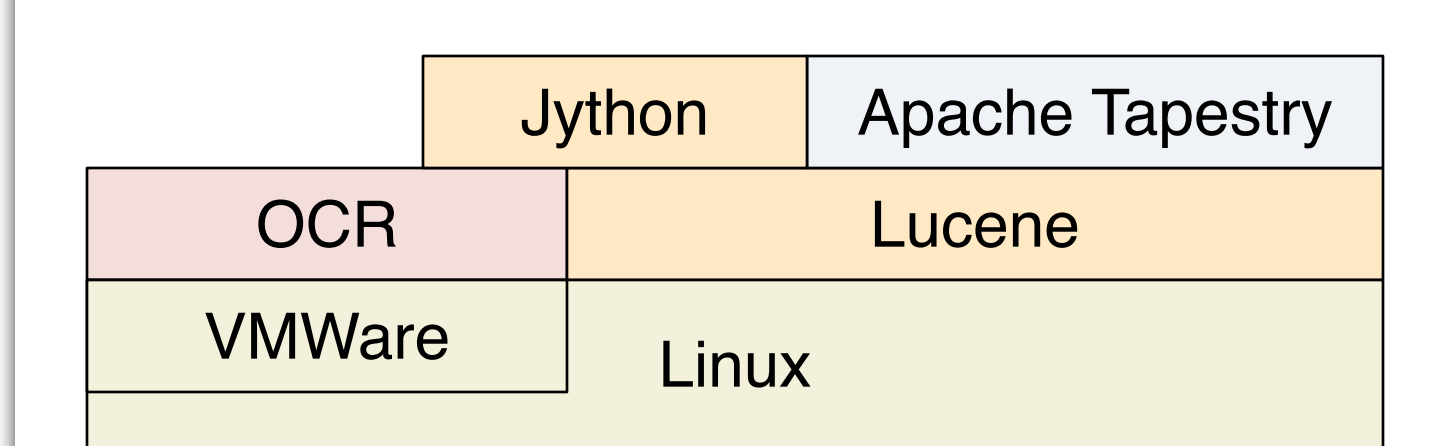| | Query | Search Target Domain | Graph | Gel | Microscopy | Diagram | List | Misc. | Total | | Relevant | | Precision |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | diet AND insulin | Caption | 17 | 0 | 0 | 0 | 0 | 2 | 19 | | 19 | | 100% |
| | | Caption & Image Text (HR) | 17 | 1 | 0 | 6 | 0 | 3 | 27 | | 25 | | 93% |
| | | Δ | 0 | 1 | 0 | 6 | 0 | 8 | 42% | 6 | 32% | | 75% |
| 2 | apoptosis AND p53 | Caption | 11 | 1 | 5 | 11 | 0 | 14 | 42 | | 42 | | 100% |
| | | Caption & Image Text (HR) | 12 | 1 | 5 | 18 | 4 | 15 | 55 | | 54 | | 98% |
| | | Δ | 1 | 0 | 0 | 7 | 4 | 1 | 13 | 31% | 12 | 29% | 92% |
| 3 | miR* AND brain AND heart | Caption | 1 | 1 | 0 | 0 | 1 | 2 | 5 | | 4 | | 80% |
| | | Caption & Image Text (HR) | 1 | 3 | 0 | 6 | 3 | 0 | 13 | | 11 | | 85% |
| | | Δ | 0 | 2 | 0 | 0 | 5 | 1 | 8 | 160% | 7 | 175% | 88% |



**Figures with short or insufficient captions are often missed by caption-based search. A combination of captions and image text offers the greatest range of searchable text. This is especially true in diagrams, lists, and heatmaps.**

## Evaluation



- plain-uncorrected
- plain-corrected
- layout-uncorrected
- layout-corrected
- corpus-plain-corrected
- corpus-layout-corrected
- high-precision
- high-recall

**Precision, Recall, and F-Rate for 8 different techniques on images with captions containing "survival" and "apoptosis". Scores are for words that are N or more characters in length. Corpus correction is less accurate on smaller words.**

## Technology

| Jython | Apache Tapestry |
|---|---|
| OCR | Lucene |
| VMWare | Linux |

## Conclusion

There are several pre- and post-processing techniques that improve OCR-based text extraction. Combinations of image layout analysis (Lienhart and Wernicke, 2002; Wu et al., 1999) and context-based correction (Kukich, 1992; Ringlstetter et al., 2007) are most beneficial. Our high recall option provides an excellent basis for text indexing and search, while our high precision option works well for more general image text extraction.

## References

M. Hearst, A. Divoli, H. Guturu, A. Ksikes, P. Nakov, M. Wooldridge, and J. Ye. 2007. Biotext search engine: beyond abstract search. Bioinformatics, 23(16):2196–2197.

K. Kukich. 1992. Technique for automatically correcting words in text. ACM Computing Surveys, 24(4): 377–439.

R. Lienhart and A. Wernicke. 2002. Localizing and segmenting text in images and videos. IEEE Transactions on Circuits and Systems for Video Technology, 12(4):256–268, Apr.

C. Ringlstetter, K. Schulz, and S. Mihov. 2007. Adaptive text correction with web-crawled domain-dependent dictionaries. ACM Transactions on Speech and Language Processing, 4(4):1–36.

V. Wu, R. Manmatha, and E.M. Riseman. 1999. Textfinder: an automatic system to detect and recognize text in images. IEEE Transactions on Pattern Analysis and Machine Intelligence, 21(11):1224–1229, Nov.

S. Xu, J. McCusker, and M. Krauthammer. 2008. Yale Image finder (YIF): a new search engine for retrieving biomedical images. Bioinformatics, Advance Access July 9, 2008.