# DOCUMENT IMAGE COMPRESSION

*by Dave Tompkins*
MASc. Candidate -- davet@ece.ubc.ca
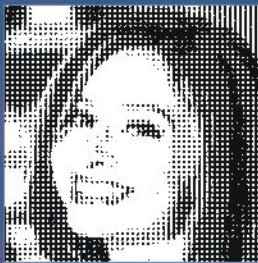Supervisor -- Dr. Faouzi Kossentini

## Our Project:

We are researching and developing advanced strategies for compressing image documents. Applications for our research include fax machines, archival systems, and Internet document systems. The framework for our research is the JBIG-2 standard.

**JBIG-2:** JBIG - the **J**oint **B**i-level **I**mage Expert **G**roup (bi-level means black & white) is an international committee comprised of industry and academic members, who have developed the standard. The standard only defines the framework and core technologies, leaving many of the compression strategies up to the people implementing the standard.

❊ *UBC is an active member of the JBIG committee, and some elements of the JBIG-2 standard have resulted from our research. We are also the first group to implement a full-featured JBIG-2 coder, and the public-domain software was written by us.*

Compare our results with some of the other available formats. If you wanted to send someone a fax by e-mail, which format would you use?

| Format | Ratio |
|--------|-------|
| BMP | 1:1 |
| GIF | 4.1:1 |
| ZIP | 5.4:1 |
| FAX | 7.8:1 |
| UBC | 14.9:1 |

| Format | Ratio |
|--------|-------|
| BMP | 1:1 |
| FAX | 1:1 |
| GIF | 2.4:1 |
| ZIP | 3.5:1 |
| UBC | 7.4:1 |

## Our Strategies:

**① SEGMENTATION**

Text, graphics, and pictures all compress differently. By *segmenting* the image, we can use a different algorithm for each type.

❊ *We have developed advanced techniques to quickly separate documents into text, line-art and pictures to improve both speed and compression.*

**② DICTIONARIES**

In any language, certain patterns of letters (*such as the word "the"*) will occur more often than others. We can exploit this to achieve higher compression. This is called a *dictionary* strategy.

Impressive Compression

Impressive Compression

❊ *We are researching advanced dictionary techniques, including language-independant dictionaries that span across several pages of a document.*

**③ REFINEMENT**

When documents are scanned or faxed, small errors are introduced so that identical letters are no longer identical. With *refinement* coding, we can adapt, so that compression increases with each subsequent letter.

e e e e

compression improves with each occurrence

❊ *We are developing pattern recognition techniques that will determine when refinement coding should be used, and how it should be implemented.*

**④ FILTERING**

By *filtering* the document and making subtle changes that are hardly noticeable to the human eye, we can significantly improve the compression performance.

rs successifs.
t pas entrepr

rs successifs.
t pas entrepr

subtle changes can improve compression

❊ *We are exploring numerous filtering strategies that achieve extremely high compression rates, without significantly distorting the image.*