

# Survey of JPEG compression history analysis

Andrew B. Lewis



Computer Laboratory

Topics in security – forensic signal analysis

## References

- ▶ Neelamani et al.:  
JPEG compression history estimation for color images, IEEE Transactions on Image Processing 15(6), 2006
- ▶ Hany Farid:  
Exposing digital forgeries from JPEG ghosts, IEEE Transactions on Information Forensics and Security 4(1), 2009
- ▶ Andrew B. Lewis, Markus G. Kuhn:  
Exact JPEG recompression<sup>1</sup>, to appear in SPIE Electronic Imaging: Visual Information Processing and Communication, 2010

---

<sup>1</sup>Draft at <http://www.cl.cam.ac.uk/~abl26/spie10-recomp-full.pdf>

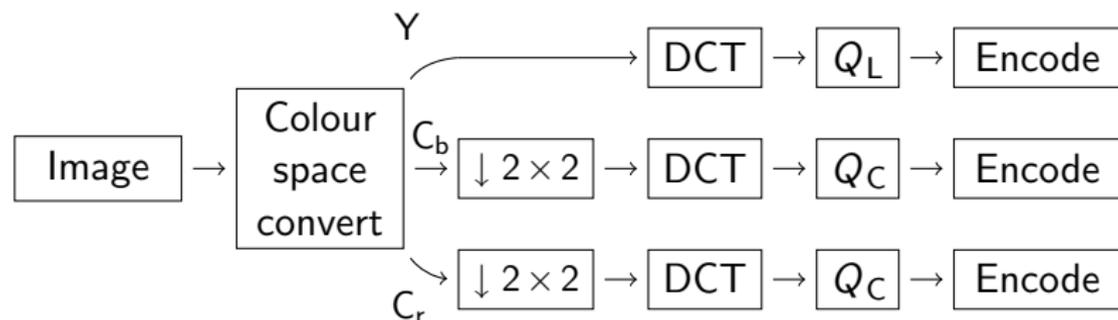
# Outline

- ▶ Revision of JPEG compression/decompression algorithm
- ▶ Probability theory for parameter estimation
- ▶ JPEG compression history estimation for color images
- ▶ Exact JPEG recompression

# The JPEG algorithm

Parameters:

- ▶ Quantization:  $Q_L$  (luma),  $Q_C$  (chroma)
- ▶ Sub-sampling:  $1 \times 1$  (luma),  $2 \times 2$ ,  $2 \times 2$  (chroma)  
also known as 4:2:0 sub-sampling
- ▶ Colour space:  $Y C_b C_r$



# JPEG compression history estimation

- ▶ Probabilistically estimate settings used in the previous compression step
- ▶ Input: raw image in colour space  $F$
- ▶ Output: compressed representation colour space  $G^*$ , sub-sampling scheme  $S^*$  and quantization tables  $Q^*$

$$\begin{aligned}\{G^*, S^*, Q^*\} &= \arg \max_{G, S, Q} P(\text{Image}, G, S, Q) \\ &= \arg \max_{G, S, Q} P(\text{Image} | G, S, Q) P(G, S, Q)\end{aligned}$$

## Terminology of inverse probability

Unknown parameters  $\theta$ , data  $D$ , assumptions  $\mathcal{H}$

$$P(\theta|D, \mathcal{H}) = \frac{P(D|\theta, \mathcal{H})P(\theta|\mathcal{H})}{P(D|\mathcal{H})}$$

$$\text{posterior} = \frac{\text{likelihood} \times \text{prior}}{\text{evidence}}$$

The quantity of  $P(D|\theta, \mathcal{H})$  is a function of both  $D$  and  $\theta$ . For fixed  $\theta$  it defines a probability over  $D$ . For fixed  $D$  it defines the likelihood of  $\theta$ .<sup>2</sup>

---

<sup>2</sup>More in David J. C. MacKay: Information Theory, Inference and Learning Algorithms, Cambridge University Press.

## Maximum a posteriori estimation

We wish to estimate  $\theta$  on the basis of data  $D$ . The maximum likelihood (ML) estimate of the parameters from the data is

$$\hat{\theta}_{\text{ML}}(D) = \arg \max_{\theta} P(D|\theta)$$

and maximum a posteriori (MAP) estimate is

$$\hat{\theta}_{\text{MAP}}(D) = \arg \max_{\theta} P(D|\theta)P(\theta)$$

## Expectation maximization

ML estimate requires the marginal likelihood, if we have hidden variables.

$$\hat{\theta}_{\text{ML}}(D) = \arg \max_{\theta} P(D|\theta)$$

$$P(D|\theta) = \sum_Z P(D|z, \theta)P(z|\theta)$$

Evaluating the sum is sometimes computationally infeasible. The expectation-maximization algorithm can be used instead.

**Expectation:**  $Q(\theta|\theta^{(t)}) = E_{Z|x, \theta^{(t)}}[\log L(\theta; x, Z)]$

**Maximization:**  $\theta^{t+1} = \arg \max_{\theta} Q(\theta|\theta^{(t)})$

No guarantees of convergence

## Interpolation characterisation as an expectation maximization problem

**Expectation:**  $Q(\theta|\theta^{(t)}) = E_{Z|x,\theta^{(t)}}[\log L(\theta; x, Z)]$

**Maximization:**  $\theta^{t+1} = \arg \max_{\theta} Q(\theta|\theta^{(t)})$

- ▶  $x$ : the observed image samples  $f(x, y)$
- ▶  $\theta$ : the interpolation kernel  $\vec{\alpha}$  and variance  $\sigma^2$
- ▶  $Z$ : the p-map, an array of probabilities with the same dimensions as the image, where each probability indicates  $P(f(x, y) \in M_1)$

## Compression history estimation as MAP problem

$$\hat{\theta}_{\text{MAP}}(D) = \arg \max_{\theta} P(D | \theta) P(\theta)$$

$$\{G^*, S^*, Q^*\} = \arg \max_{G, S, Q} P(\text{Image} | G, S, Q) P(G, S, Q)$$

## Estimating quantization tables $Q^*$ (1)

- ▶ Small set of possible values for  $G$  and  $S$ .
- ▶ C'space to  $G^*$ , sub-sample with  $S^*$ , then forward DCT gives a near-periodic distribution of coefficients  $\Omega_{G,S}$  over image.
- ▶  $G, S, Q$  and  $\tilde{X}_{G,S} \in \Omega_{G,S}$  independent

$$\{G^*, S^*, Q^*\} = \arg \max_{G,S,Q} P(\Omega_{G,S}|G, S, Q)P(G)P(S)P(Q)$$

$$\{G^*, S^*, Q^*\} = \arg \max_{G,S,Q} \prod_{\tilde{X}_{G,S} \in \Omega_{G,S}} P(\tilde{X}_{G,S}|G, S, Q)P(G)P(S)P(Q)$$

## Estimating quantization tables $Q^*$ (2)

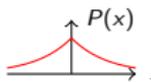
Since the decompressor's dequantization,

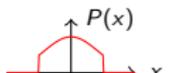
- ▶ DCT coefficients  $\bar{X}_q$  were IDCT'ed;
- ▶ the results were up-sampled if necessary; and
- ▶ the image was converted to the RGB colour space.
- ▶ We received the image, and applied a forward colour space conversion;
- ▶ we downsampled the planes, if appropriate; and
- ▶ we applied the forward DCT to get  $\tilde{X}$ .

Rounding errors accumulate during every stage of this process. Therefore, we model the DCT coefficient values

$$\tilde{X} = \bar{X}_q + \Gamma$$

where original DCT coefficients  $\bar{X}_q$  are modelled by a sampled

zero-mean Laplace distribution  and rounding error  $\Gamma$  is

drawn from a truncated normal distribution .

## Estimating quantization tables $Q^*$ (3)

The Laplace distribution for DCT coefficient values has scale parameter  $\lambda$ , determined from the observations. The probability distribution of coefficient values after quantization/dequantization with factor  $q$  is

$$P(\bar{X}_q = t | q \in \mathbb{Z}^+) = \sum_{k \in \mathbb{Z}} \delta(t - kq) \cdot \int_{(k-0.5)q}^{(k+0.5)q} \frac{\lambda}{2} \exp(-\lambda|\tau|) d\tau$$

Rounding errors are independent, so we convolve this distribution with our error term's distribution and normalize:

$$P(\tilde{X} = t | q) \propto \int P(\bar{X}_q = \tau | q) P(\Gamma = t - \tau) d\tau$$

## Estimating quantization tables $Q^*$ (4)

If we assume a uniform prior  $P(q)$  for quantization factors, we can now find the most likely value for a particular quantization factor  $q = Q_{i,j}^*$  by maximizing  $P(\tilde{X}|q)$ . This is repeated for each quantization factor  $Q_{i,j}^*$  for  $(i,j) \in \{(0,0), \dots, (7,7)\}$ :

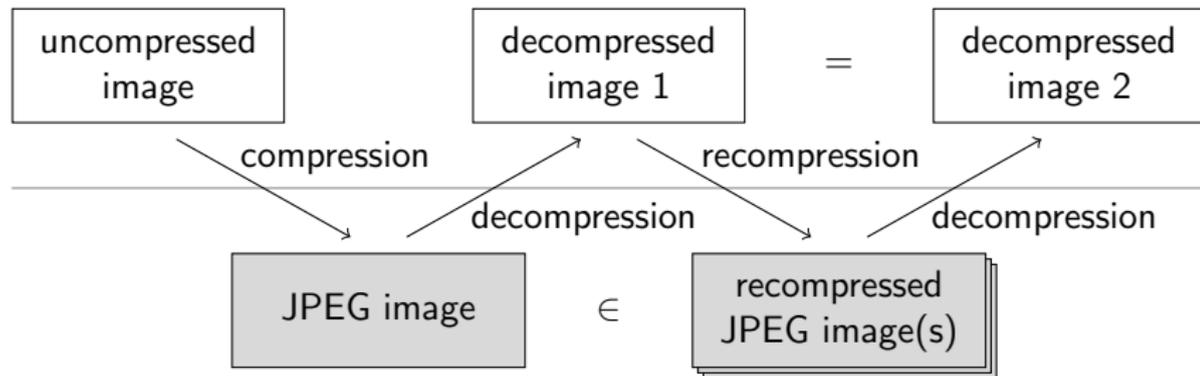
$$q^* = \arg \max_{q \in \mathbb{Z}^+} \left( \prod_{\tilde{X} \in \Omega} P(\tilde{X}|q) \right)$$

# Evaluation

- ▶ Allows for arbitrary colour space conversion, sub-sampling and quantization table parameters
- ▶ Not implementation specific
- ▶ Partially recovered quantization tables could be used to determine quality factor  
Quality factors  $q \in \{1, \dots, 100\}$  map onto quantization tables in many compressor implementations
- ▶ Statistical rather than exact
- ▶ Can't recover bitstream
- ▶ Tikhonov deconvolution filter introduces errors
- ▶ Errors in the quantization table when most DCT coefficients at a particular frequency are zero
- ▶ No data for the quantization table when all are zero (low quality factors)

## Exact recompression

- ▶ Can we recover the original bitstream (including the quantization tables) when given the result of decompression?
- ▶ Due to rounding and mismatch between the compressor and decompressor operations, simply invoking the compressor with the same parameters will not work.
- ▶ Can we provide a guarantee that if the provided image was produced by a particular JPEG decompressor, we will recover that bitstream?



## Applications of exact recompression

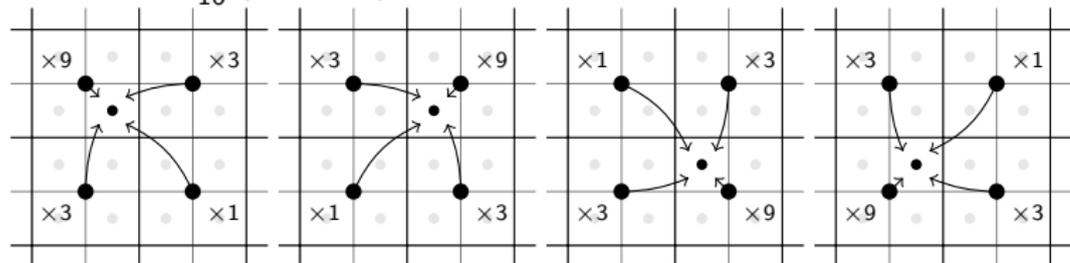
- ▶ The input to JPEG compressors is often previously compressed image data. Detecting this and recompressing exactly will reduce the information loss from recompression.
- ▶ Hinder forensic analysis – double compression detection, JPEG ‘ghosts’, . . .
- ▶ Detect tampered regions in an uncompressed image, when the background was output by a JPEG decompressor
- ▶ Some copy-protection schemes rely on the fact that copies will be recompressed, lowering the quality.

## Information loss in the decompressor

- ▶ Model sources of uncertainty (rounding, range limiting) when inverting operations in the decompressor using intervals of integers. We store intervals for all intermediate values in the decompressor.
- ▶ Because we are modelling the exact computations, exact recompressors are implementation-specific.

## Chroma down-sampling example (1)

IJG chroma upsampling filter weights contributions from neighbouring samples by  $\frac{1}{16}$  (1, 3, 3, 9) in order of increasing proximity.



$$v_{x,y} = \left[ \frac{1}{16} (8 + \alpha \cdot w_{i-1,j-1} + \beta \cdot w_{i,j-1} + \gamma \cdot w_{i-1,j} + \delta \cdot w_{i,j}) \right],$$

with weights

$$(\alpha, \beta, \gamma, \delta) = \begin{cases} (1, 3, 3, 9) & x = 2i, y = 2j \\ (3, 1, 9, 3) & x = 2i - 1, y = 2j \\ (3, 9, 1, 3) & x = 2i, y = 2j - 1 \\ (9, 3, 3, 1) & x = 2i - 1, y = 2j - 1. \end{cases}$$

## Chroma down-sampling example (2)

We need to solve for the down-sampled weights  $w_{i,j}$ :

$$v_{x,y} = \left[ \frac{1}{16} (8 + \alpha \cdot w_{i-1,j-1} + \beta \cdot w_{i,j-1} + \gamma \cdot w_{i-1,j} + \delta \cdot w_{i,j}) \right]$$

Interval arithmetic rules give

$$\bar{w}_{i,j} = \left[ \left[ \frac{1}{\delta} (v_{x,y} \perp \times 16 - (8 + \alpha \cdot \bar{w}_{i-1,j-1} + \beta \cdot \bar{w}_{i,j-1} + \gamma \cdot \bar{w}_{i-1,j})) \right], \right. \\ \left. \left[ \frac{1}{\delta} (v_{x,y} \top \times 16 + 15 - (\alpha \cdot \bar{w}_{i-1,j-1} + \beta \cdot \bar{w}_{i,j-1} + \gamma \cdot \bar{w}_{i-1,j})) \right] \right]$$

## Chroma down-sampling example (3)

$k \leftarrow 0$

$\bar{w}_{x,y}^0 \leftarrow [0, 255]$  at all positions  $-1 \leq x \leq \frac{w}{2}, -1 \leq y \leq \frac{h}{2}$

**repeat**

$k \leftarrow k + 1$

change scan order of  $(x, y)$  ( $\begin{smallmatrix} \rightarrow \\ \vdots \\ \rightarrow \end{smallmatrix}$ ,  $\begin{smallmatrix} \leftarrow \\ \vdots \\ \leftarrow \end{smallmatrix}$ ,  $\begin{smallmatrix} \rightarrow \\ \vdots \\ \leftarrow \end{smallmatrix}$ ,  $\begin{smallmatrix} \leftarrow \\ \vdots \\ \rightarrow \end{smallmatrix}$ )

**for** each sample position  $(x, y)$  in the upsampled plane **do**

**for**  $(i', j') \in \{(i-1, j-1), (i, j-1), (i-1, j), (i, j)\}$  **do**

$\bar{w}'_{i',j'} \leftarrow \bigcup_{s \in \check{v}_{x,y}^c} \bar{a} :$

Equation satisfied with  $\bar{a}$  for  $w_{i',j'}$ ,  $s$  for  $v_{x,y}$

and current estimates  $\bar{w}_{x,y}$  for other  $w$  values.

$\bar{w}'_{i',j'} \leftarrow \bar{w}'_{i',j'} \cap \bar{w}_{i',j'}^{k-1}$

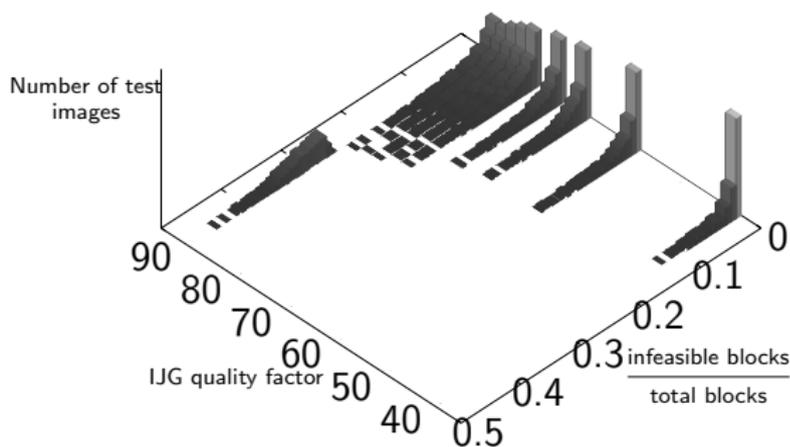
**end for**

**end for**

**until**  $\bar{w}^k = \bar{w}^{k-1}$

$\bar{w}_{x,y}^c \leftarrow \bar{w}^k$

# Performance



**Figure:** Recompression performance for a dataset of uncompressed images from the UCID. 1338 colour images were compressed and decompressed at quality factors  $q \in \{40, 60, 70, 80, 82, 84, 86, 88, 90\}$ , then recompressed. The proportion of blocks at each quality factor which were not possible to recompress due to an infeasible search size is shown.