

# Cloud-scale Image Compression Through Content Deduplication

David Perra  
perra@cs.unc.edu

Jan-Michael Frahm  
jmf@cs.unc.edu

Department of Computer Science,  
University of North Carolina,  
Chapel Hill, NC

Department of Computer Science,  
University of North Carolina,  
Chapel Hill, NC

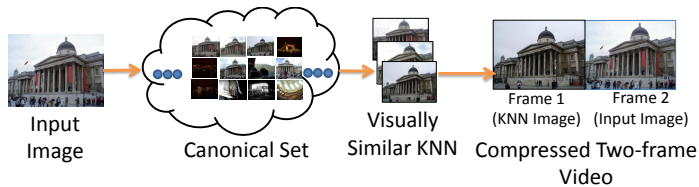


Figure 1: A high-level overview of our proposed technique. A canonical set, consisting of potentially millions of images of varying subjects, can be leveraged to find a set of images which are visually similar to a query image via  $k$ -nearest neighbors search. Those visually similar images are each used to compress the query image by use of a video codec such as H.265. The two-frame video which yields the best compression to quality ratio will be selected as the output of the pipeline, and the bits associated with the input image will be stored.

Modern large-scale photo services such as Google Drive, Microsoft OneDrive, Facebook, and Flickr are currently tasked with storing and serving unprecedented quantities of photo data. While most photo services still utilize jpeg compression to store photos, more elegant compression schemes will need to evolve to combat the storage costs associated with the exponential increase in data. To satisfy this need, two classes of solutions have been established: representative signal techniques [1, 6], and visual clustering techniques [3, 5, 8]. Representative signal techniques work by finding a common low-frequency signal within a set of images. The technique presented in this paper falls into the second class of techniques, which focuses upon sharing and reusing pixel data between multiple images by modelling the relationship between these images as a directional graph. Paths through this directional graph define image pseudosequences, or directionally related subsets of images which describe the visually shortest path between various images in an image set [5, 7, 8]. These pseudosequences can then be used for compression via image reconstruction or compression via traditional video codecs, such as H.265.

The primary shortcomings for state-of-the-art visual clustering techniques stem from a lack of scalability. Finding appropriate image pseudosequences becomes increasingly more difficult as an image set grows. This is because all-pairs comparisons must be performed between the images to find an optimal graph. Additionally, longer pseudosequences tend to result from larger image sets. Longer pseudosequences cause image compression and decompression to take longer, leading to a decrease in performance with an increase in dataset size.

In this paper, we present an efficient cloud-scale digital image compression scheme which overcomes the scalability issues found in the state-of-the-art techniques. Unlike current state-of-the-art systems, our image compression technique takes full advantage of redundant image data in the cloud by independently compressing each newly uploaded image with its GIST nearest neighbor taken from a canonical set of uncompressed images. This allows for fast identification of a size-restricted pseudosequence. We then leverage state-of-the-art video compression techniques, namely H.265, in order to efficiently reuse image content which is already stored server-side.

Previous state-of-the-art techniques used only the image data found within a particular image set to compress the entire set [5, 7, 8]. Our technique, on the other hand, avoids this through the use of the canonical set of images. Our key insight is that many photographs uploaded to the cloud are highly likely to have similar pixel patches, especially near landmarks and other commonly geotagged areas – even within the home of the user. Thus, we assume that the canonical set is a randomly selected, finite set of photos that is composed of tens or hundreds of millions of images depicting commonly photographed subjects and structures. Con-

structing such a set can be done, for example, by randomly sampling all photos currently stored in the cloud. Alternatively, techniques like Frahm *et al.* [2] and Raguram *et al.* [4] can be used to construct such a canonical set through iconic scene graphs. This process should naturally yield many views of popular subjects as more photos of those subjects are uploaded to the cloud. A sufficiently large canonical set contains enough photos to have a visually similar image for a large majority of photos uploaded in the future. Similarly, we foresee complementing the general canonical set with a user-specific canonical set if desired. Once an ideal canonical set is constructed, it can be used as a generic dataset for compressing any photo collection.

The implementation of our method is described in our paper, and extensive experiments are conducted. Experimental results demonstrate that our algorithm produces competitive image compression rates while reducing the computational effort by at least an order of magnitude in comparison to competing techniques, all while providing the necessary scalability for use in cloud-scale applications.

- [1] Samy Ait-Aoudia and Abdelhalim Gabis. A comparison of set redundancy compression techniques. *EURASIP J. Appl. Signal Process.*, 2006:216–216, January 2006.
- [2] Jan-Michael Frahm, Pierre Fite-Georgel, David Gallup, Tim Johnson, Rahul Raguram, Changchang Wu, Yi-Hung Jen, Enrique Dunn, Brian Clipp, Svetlana Lazebnik, and Marc Pollefeys. Building rome on a cloudless day. In Kostas Daniilidis, Petros Maragos, and Nikos Paragios, editors, *ECCV*, volume 6314 of *Lecture Notes in Computer Science*, pages 368–381. Springer Berlin Heidelberg, 2010.
- [3] Yang Lu, Tien-Tsin Wong, and Pheng-Ann Heng. Digital photo similarity analysis in frequency domain and photo album compression. In *Proceedings of the 3rd International Conference on Mobile and Ubiquitous Multimedia*, MUM '04, pages 237–244, New York, NY, USA, 2004. ACM.
- [4] Rahul Raguram, Changchang Wu, Jan-Michael Frahm, and Svetlana Lazebnik. Modeling and recognition of landmark image collections using iconic scene graphs. *Int. J. Comput. Vision*, 95(3):213–239, December 2011.
- [5] Zhongbo Shi, Xiaoyan Sun, and Feng Wu. Photo album compression for cloud storage using local features. *Emerging and Selected Topics in Circuits and Systems, IEEE Journal on*, 4(1):17–28, March 2014.
- [6] Chi-Ho Yeung, O.C. Au, Ketan Tang, Zhiding Yu, Enming Luo, Yannan Wu, and Shing-Fat Tu. Compressing similar image sets using low frequency template. In *Multimedia and Expo (ICME), 2011 IEEE International Conference on*, pages 1–6, July 2011.
- [7] Huanjing Yue, Xiaoyan Sun, Jingyu Yang, and Feng Wu. Cloud-based image coding for mobile devices 2014; toward thousands to one compression. *Multimedia, IEEE Transactions on*, 15(4):845–857, June 2013.
- [8] Ruobing Zou, O.C. Au, Guyue Zhou, Wei Dai, Wei Hu, and Pengfei Wan. Personal photo album compression and management. In *Circuits and Systems (ISCAS), 2013 IEEE International Symposium on*, pages 1428–1431, May 2013.