# Herding Cats

**John R. Smith**
*IBM T. J. Watson Research Center*

The more things change, the more they stay the same. Over the last decade, digital video has gone from nowhere to everywhere. Today video makes up more than half of mobile and peak Internet traffic, and more video is being captured in a second than anyone can view in a lifetime. According to YouTube, one of today's most popular video sites, 100 hours of video are uploaded to its site every minute and more than 6 billion hours of video are watched each month (see www.youtube.com/yt/press/statistics.html).

But therein lies the problem. Searching for video is still difficult, and techniques for browsing and visualizing video content are primitive at best. The predominant method for video search today is based on user tags and limited metadata. If the search corpus is large enough and meta information is available, video search can often find some relevant videos. However, there are problems with both precision and recall in video search today. This means greater emphasis is needed on interactive tools that allow fast and effective triage of video search results.

As an illustrative example, the video searches in Table 1 were performed on YouTube. The table shows the number of results found for each of seven video searches. Although the sheer number of cat videos that can be searched on YouTube is truly amazing, there are problems, as the table shows. Searching for "kittens" provides more than 3 million results, which is way beyond what any user can scan. Attempts to make the search more specific by searching for "calico kittens playing" or "calico kittens playing outdoors with dog and children" reduce the number of results significantly. Still, the result size is still in the tens of thousands.

Finally, searching for "calico kittens playing outdoors with dog and children and cat toy" reduces the number of results to 2,570. However, these results look sketchy and more precise queries seem to reduce both precision and

*Table 1. Video searches on YouTube.*

| Search | Number of results |
|---|---|
| Kittens | 3,340,000 |
| Calico kittens | 55,600 |
| Calico kittens playing | 42,400 |
| Calico kittens playing outdoors | 34,100 |
| Calico kittens playing outdoors with dog | 37,000 |
| Calico kittens playing outdoors with dog and children | 10,400 |
| Calico kittens playing outdoors with dog and children and cat toy | 2,570 |

recall. If users are looking for specific content, there's really nothing they can do to find the relevant videos except click and watch each one.

This is a situation where compact visual summaries are needed to make video triage easier. The most basic video summary is a thumbnail, which is what YouTube provides. Thumbnail images have been around since the beginning of the Web and are commonly used for image search. Providing a single frame thumbnail is obviously an impoverished visual summary for video given it's a temporal medium. Some sites go further and provide an animated thumbnail that displays a sequence of frames from the video. More advanced techniques base the animated sequence on automatically detected keyframes from the video. Although animations can provide a better summary, they display their information over time, which can slow down triage. Alternatively, techniques like video filmstrips and mosaics can display the keyframe information in 2D space, making it easier to see more of the video content at a glance. However, if videos are long, these representations can be dense or take up too much screen real estate.

There may not be a single, best video summary for all situations. Depending on the context of the user's search and other factors such as screen size, bandwidth, and number of search results, there should be different trade-offs in space, time, and content for making effective and efficient video summaries.

Video has indeed come a long way. But video search has not. Video summarization and search go hand in hand, and more effective techniques are needed for both. Until then, finding relevant videos will remain akin to herding cats. **MM**

**John R. Smith** is a senior manager of Intelligent Information Management at IBM T.J. Watson Research Center. Contact him at jsmith@us.ibm.com.