
Human Computation in Online Video Storytelling

Philo D. I. van Kemenade

Submitted in partial fulfillment of the requirement of the degree of
MASTER OF SCIENCE IN COGNITIVE COMPUTING

Goldsmiths College
University of London

Supervisor

Dr. Marian Ursu

Department of Computing
Goldsmiths, University of London
New Cross, London SE14 6NW

September 16, 2012

Abstract

Digital video retrieval, filtering and reconfiguration are difficult tasks to solve using current computational techniques. An important cause of this difficulty is the semantic gap between a visual representation and the meaning we address to it. A solution commonly sought in AI research is to reduce the gap by visual analysis and the linking thereof to previously established symbolic representations of semantical concepts. These methods often perform poorly on unpredictable content found in the large video libraries of user generated content that account for much of the global internet traffic these days. A second way of hunting down meaning in visual content is to step over the gap altogether and ask people directly for a meaningful interpretation one wishes to acquire for an item of content. By accessing many people's interpretations in small bite-sized tasks, collectively grounded annotations can be established. This form of accessing human computational power has seen a major increase in attention and application, for a large part because of the increased connectivity of individuals to the web. This thesis investigates how tasks involving meaningful interpretation of video content can benefit from the use of human computation. In order to test the validity of these approaches 'wePorter' is developed, a system with the purpose of finding local intervals of interest within videos in a larger set of topically related content. We also investigate how such a system can be used for reconfiguration of content into new and informative stories. [concluding]

Contents

1	Introduction	4
2	The Quest for Meaning in Video	5
2.1	Forms of Meaning in Film and Video	5
2.1.1	Meaning in Visuals	5
2.1.2	Meaning in Concept	5
2.1.3	Meaning in Structure	5
2.1.4	Meaning in Annotations	5
2.2	Computational Undertakings of the Quest	5
2.2.1	The Semantic Gap	5
2.2.2	Steps towards Meaning: An Overview	5
2.2.3	Indexing to enable search	5
3	Human Computation towards Visual Meaning	6
3.0.4	The web as platform for creation	7
3.1	Deriving Meaning from Video Via Human Factors	7
3.2	Collaborative Filtering	7
3.3	A Characterisation of Human Computation Systems	7
3.3.1	Purpose	7
3.3.2	Motivation	7
3.3.3	Task	7
4	Interactive Storytelling: From database to data-based	8
4.1	Symbolic approaches	8
4.2	Statistic approaches	8
4.3	Remixing	8
4.3.1	Taking the remix online	8
5	Human Computed Stories in wePorter	9
5.1	Introduction	9
5.2	User Generated Video Content	9
5.3	The Purpose	10
5.3.1	Different Goals	10
5.3.2	Serving the Purpose of Unarticulated Want	12
5.3.3	Storytelling as Structured Recommendation	12
5.4	The Motivation	13
5.4.1	Information Provision through Online Video	13
5.5	The Task	14
5.5.1	Design Considerations	14

5.5.2	The Interface	17
5.5.3	Forced Feedback	20
5.6	Analysis of the Weporter System	20
5.6.1	Landscapes of Interest	20
5.6.2	The death of the author	20
5.7	Parallel Play	20
5.8	Implementation	20
6	Evaluation	22
7	Discussion	23
8	Future Directions	24
9	Conclusions	25
	Bibliography	26

Chapter 1

Introduction

Chapter 2

The Quest for Meaning in Video

2.1 Forms of Meaning in Film and Video

2.1.1 Meaning in Visuals

[Video structuring, indexing and retrieval based on global motion wavelet coefficients][2]

2.1.2 Meaning in Concept

2.1.3 Meaning in Structure

[cite bordwell & Thompson: analysis of context dependency]

2.1.4 Meaning in Annotations

[Telop-on-demand: Video structuring and retrieval based on text recognition][13] [Addressing the Challenge of Visual Information Access from Digital Image and Video Libraries][6]

2.2 Computational Undertakings of the Quest

2.2.1 The Semantic Gap

2.2.2 Steps towards Meaning: An Overview

Schematized summary of different steps: indexing automatic metadata annotating Human-Driven Labeling Machine-Driven Labeling multimodal feature fusion concept based content based [Relevance feedback: A power tool for interactive content-based image retrieval][20] [The Relative Effectiveness of Concept-based Versus Content-based Video Retrieval][26] [collaborative filtering] http://en.wikipedia.org/wiki/Collaborative_filtering

2.2.3 Indexing to enable search

Because visual data on it's own provides little machine-readable handles to search and find, repositories of multimedia content need to be index to enable search. Within the task of video indexing several approaches are taken

Chapter 3

Human Computation towards Visual Meaning

[Define GWAP]

In a recent survey paper, Quinn and Bederson present a taxonomy of Human computation systems[17]. They sketch out the trend of this new method for intelligent problem solving by the increase of academic papers featuring the term ‘human computation’ and its relative ‘crowd-sourcing’. They summarise the myriad of definitions given in recent works by several different authors in two key points:

- “The problems fit the general paradigm of computation, and as such might someday be solvable by computers.”
- “The human participation is directed by the computational system of process”

The first point introduces an interesting question whether storytelling is a computable process. In 1950, Alan Turing envisioned in his seminal paper ‘Computing Machinery and Intelligence’ that a computer program would be able to successfully play a game now known as the Turing Test. His work also mentions that

“[t]he idea behind digital computers may be explained by saying that these machines are intended to carry out any operations which could be done by a human computer” [22]

The notion of ‘human computer’ benefits from some contextualisation, as in the last few decades we’ve become unaccustomed to the term. Human computers were not uncommon in the time of Turing and before that from the 18th century, when ‘computer’ was used to signify ‘one who computes’[10]. People bearing the function title were involved in the execution of calculations produced by strictly following mathematical theories. The activity that these computers were involved in was a process of rote, not requiring any human creativity. While working on the design for the first ever mechanical computer, Charles Babbage called it “mental labour”[1, Ch. 20].

Quinn and Bederson further present a classification along six dimensions they see as the most salient distinguishing factors:

[table of dimensions?]

3.0.4 The web as platform for creation

Many media scholars have written about the role of the web [refs New Media Reader]. Important trend of the web as platform of creation. In terms of video creation for example, the last few years have seen the development of online video editing tools and environments such as popcorn.js, WeVideo and Kaltura.

3.1 Deriving Meaning from Video Via Human Factors

[Personalized online document, image and video recommendation via commodity eye-tracking][24]

[VideoReach: an online video recommendation system][15]

3.2 Collaborative Filtering

The idea of using user's past interactions within a system hosting digital content for the filtering of items that might be of interest is not a new one and usually goes by the name of collaborative filtering. Collaborative filtering can generally take two forms: User-based, Item-based

Information filtering agents and collaborative filtering both attempt to alleviate information overload by identifying which items a user will find worthwhile. I

already happening at YouTube

3.3 A Characterisation of Human Computation Systems

3.3.1 Purpose

3.3.2 Motivation

3.3.3 Task

Chapter 4

Interactive Storytelling: From database to data-based

4.1 Symbolic approaches

[28, 18, 23]

4.2 Statistic approaches

4.3 Remixing

The reconfiguration of smaller units that carry meaning within themselves is common practice for textual media such as blogs, where it is easy to quote part of another author's writing in a new post [ref].

It needs to be said that some reconfiguration of videos is taking place, but even though it often concerns content that was originally sourced online, much of the creative act of remixing happens offline.

4.3.1 Taking the remix online

examples like: Aaron Koblin (johnny cash project, exquisite corps, etc)

there is even word of a true remix culture.

the availability of multi-media content via the internet has meant a surge in

Chapter 5

Human Computed Stories in wePorter

[intro: sketch wePorter scenario. youtube: Burning Man 1.7 million, Burning Man 2012 ..., Burning Man this month ...]

[discuss: is wePorter human computation? is it solving a task that would otherwise be done by computers? By humans? Can it be expected to ever be done by computers? Looking at answers from the questionnaire to the question of why users decided to change their focus, reveal different reasons for their behaviour. A big part of the answers are generic visual or auditory features like the presence of bright colours, video or audio quality and sudden changes in movement of the scenery of camera. These are all aspects of a video that can be computed computationally, some day perhaps even in a realtime online environment. The question whether or not other features such as curiosity or boredom are computable or not is beyond the scope of this thesis, but participants' response hint at the possibility that their interest is at least partially explicable by computable factors.

plot aggregated reasons to change:

cinematography camera movement content based action trees, boats visual colour quality changes, movement audio music Interest, curiosity / boredom unpleasant sound visuals
]

5.1 Introduction

This section describes the interaction design of the wePorter system, built to exemplify how human computation can be used in tasks like local video part filtering and semi-automated video reconfiguration. The wePorter system runs an interactive webpage that functions as the main source for data acquisition, presentation of results and general proof of concept. In this chapter the system is analysed along the axes of *Purpose*, *Motivation* and *Task*, that were introduced in the analysis of Human Computational systems in chapter 3. Next to these guidelines for analysis, some remarks are made about the specifics in the functioning of the system. Lastly we discuss the implementation of the web application that is central in wePorter.

5.2 User Generated Video Content

[define UGVC]

http://youtube-global.blogspot.co.uk/2009/05/zoinks-20-hours-of-video-uploaded-every_20.html
<http://youtube-global.blogspot.co.uk/2010/03/oops-pow-surprise24-hours-of-video-all.html>
<http://youtube-global.blogspot.co.uk/2010/11/great-scott-over-35-hours-of-video.html>
<http://youtube-global.blogspot.co.uk/2011/05/thanks-youtube-community-for-two-big.html>
<http://youtube-global.blogspot.co.uk/2012/01/holy-nyans-60-hours-per-minute-and-4.html>

Since the dawn of YouTube, we've been sharing the hours of video you upload every minute. In 2007 we started at six hours, then in 2010 we were at 24 hours, then 35, then 48, and now...60 hours of video every minute, an increase of more than 25 percent in the last eight months. In other words, you're uploading one hour of video to YouTube every second. Tick, tock, tick, tock that's 4 hours right there!

These astonishing figures of the amount of video that is uploaded to youtube are nothing short of mind blowing, but will most likely sound dated in a matter of years or even months. Looking at the increase of content uploaded to the video platform in past years, the growth does not seem likely to come to a halt soon [ref table]. All these videos are great for online video junkies, and are increasingly part of the online journalism landscape [19]. At the same time, all these videos being put online beg the question which ones of them to watch.

[table of youtube content uploads]

The increasing amounts of content being put online, lead to an information overload and present serious challenges in search and information retrieval tasks [ref]. There is an increased need for ways of aggregation and filtering. Both of these tasks rely heavily on an at least a shallow understanding of what is presented in these media, which, as we've seen in chapter 2, is a hard problem to solve via current computational techniques. With so much content being uploaded, how can we find our way in the already enormous ocean of online videos?

5.3 The Purpose

Searching -¿ IR With more than an hour of new content per second it is no wonder that youtube has come to be viewed as the go-to for online video, much like "the digital video repository for the Internet" that was envisioned by its founders in their first ever blog post [ref <http://youtube-global.blogspot.co.uk/2005/07/greetings-everyone-thanks-for-visiting.html>]. An important activity on video platforms like youtube is searching and much attention has been given to different methods of multimedia search and indexing [refs]. Youtube's acquisition by Google in 2006 underlines the platform's role as a video search engine.

5.3.1 Different Goals

Once items have been annotated with tags that reflect the content of a video, these indexes can, along with other meta data of the video, be used for retrieval of videos[ref] in response to textual queries. The effectiveness of such a retrieval task can vary depending on the information that is used in the search algorithm[refs] and the type of content that is searched for [11][more refs]. A third characteristic that determines the effectiveness of a

video retrieval system is the goals that users have in their usage of the system. These user goals can vary widely from more to less specific[7]. We expand on this latter point, as it forms an important context for the wePorter system.

Direct Navigation

First, a user might be drawn to a video platform by a direct link from an external website. Links might be either in the form of actual hyperlinks or playable embedded videos that are followed through to the platform. Navigations via such links form a direct mapping between a user's intention to the desired piece of content. In this case, users have a very specific reason to come and watch. Their desire, at least of knowing the contents of the video, is satisfied after the viewing. YouTube's system engineers call this way of video viewing *direct navigation*[7].

Search and Goal-oriented browse

When users have not obtained a direct link to a potentially relevant piece of content, they might still have a specific goal in mind when visiting an online video platform. Reasons to visit might be the wish to see a particular music video or to find an instance of a series by a particular producer. This goal of discovering a rather specific video is referred to as *search and goal-oriented browse*. Provided that the desired piece of content exists and the video platform has an appropriate search function in place, these 'narrow queries', will result in a result set of search results from which the user is likely to handpick the sought-after result fairly quickly. Here the user's desired result often lies within a single item of content. Perhaps a few hit and misses are required, but after a couple of clicks the required video is found.

Unarticulated Want

[first general unarticulated want, then specifically topically encapsulated browse]

A different scenario emerges altogether when a user approaches a video platform with a more broad and open ended motive. Think for example of someone who wants to know what happened at a large music festival she recently attended, or someone who couldn't make it to a large demonstration and would like to get a sense of the atmosphere. This kind of 'broad queries' returns a result set of content in which a user will probably consider many items as a successful retrieval. Furthermore, one could even say that the desired result of a user's query is spread across the multiple pieces of content. By traversing the space of different videos in the result set, users interactively construct the answers to their own queries.

This kind of navigation through a space of related content is common practice on large spaces of linked data[ref]. It has been found that YouTube's related video recommendation functionality, which recommends videos that are related to the video currently being watched, is one of the most important view sources of videos. In fact, traffic received from these recommendations is the main source of views for the majority of videos on YouTube [27]. Features derived from users' navigations such as 'click-through rate' have been used to improve content-based video recommendation [25].

Goal of a person's query in this kind of navigation is no longer defined in a single returnable item of content or even a containable set of items. Rather, the interactive pathway through the a set of interesting bits of content is what represents a user's aim. This broader exploratory goal of finding different parts of relevant content has been termed 'unarticulated want'[7].

Interactivity is generally agreed to play an important role in the task of video retrieval, as is reflected by the separate category in the annual TRECvid challenge for interactive video retrieval[21]. Several works have indicated the importance of interactivity in the task of video retrieval to filter through a set of initially returned results [8, 5, 8, 9]. While most of these systems are aimed at retrieval of clearly specified queries, exemplified by the TRECvid retrieval task, the need for interactive exploration is even more apparent for the broader oriented goal of users engaged in unarticulated want.

5.3.2 Serving the Purpose of Unarticulated Want

The answer to a user’s query now lies as much in the journey through the content as in the returned content itself. By traversing from one piece of content to the next, users construct their a sequence of concatenated pieces content. This self-constructed story is an important concept that wePorter capitalises on, as will soon become apparent.

The task at hand of recommending a larger group of interesting videos is radically different compared to the more narrow queries that could be answered by a small set of true positives in an information retrieval task. Besides the spread of the searched for result across different pieces of content, there is a second important difference that lies in the nature of the majority of UGVC.

User-contributed videos commonly consist of raw, unedited footage. In [19] Rosentiel and Mitchell report that within the collection they investigated only 39% of the news-related footage contributed by citizens was edited. It should be noted that this collection contained only the most popular videos per week and that a different distribution will be found in the complete set of news-related videos or all the videos hosted on YouTube.

Users with broad expectations will not only want to be presented with multiple relevant items from a complete repository, they are also looking for the most interesting parts within these relevant items. This issue is particular to time-based media, and especially relevant for video. Other temporal media, like audio in general and music in particular, have less of a need for segmentation because of their common usage in multimedia applications. People usually tend to listen to a song entirely and if they which to experience an album in part, constituent songs are already units on their own that can easily be reconfigured. Tag-a-tune is a game with a purpose used to acquire tags for clips of music. Although it could be employed for labelling of smaller audio subclips within songs, the games only aims at global labelling of a sounds[14].

Because of the raw, unedited nature of the majority of UGVC it is desirable to establish local recommendations that point to ‘subclips’ within a video that are of particular interest. Whereas digital music albums shared online consist of a collection of songs that can each easily be made to stand alone, video currently suffers from a less malleable identity online. Online videos are currently much like black boxes that can be played, paused, rated, commented on, tagged and shared only in its entirety. What if a piece of raw, unedited UGVC features something spectacular for ten seconds halfway along its timeline, but shows much of the same for the rest of the time? Answering this question will be the first part of the purpose of the wePorter system.

5.3.3 Storytelling as Structured Recommendation

explanation why it is also good to structure the filtered pieces of content into a new configuration. [TODO]

The ten significant seconds in a two-minute video become a needle in a haystack when an initial set of videos relating to your query includes tens to hundreds of possibly relevant

videos with lengths between some tens of seconds and a couple of minutes. The aggregation and reconfiguration of several of these ‘needles’ into a meaningful new whole is another non-trivial task. We present wePorter as a test case for new methods that address both these issues of information overload in video libraries of UGVC. More precisely, wePorter’s purpose is two-folded:

From a set of topically related unedited user-generated videos:

1. Find localised intervals of interest within each of the source videos
2. Find a meaningful structure for the reconfiguration of interesting video parts

Idea 1: towards interest based filtering via meaning

Idea 2: skip the semantic gap and directly model human attentional behaviour

many people watching: Clicking from one video to the next (choosing from a set of related videos) these inter-video links could be seen as indicators for relatedness and relevance, much like google’s page rank algorithm use links across webpages to establish a notion of the most significant site on a particular topic.

There is an important difference here though. Whereas the links used by Google’s search algorithms are embedded in machine readable hyperlinks, the path of clicking on from one video to the next is a characteristic of a person’s interaction.

differences: public, readable // private, non readable conscious choice // unconscious result of interaction Concluding can be consciously put in place by several people at large scale // dependent on real ‘human’ traffic.

5.4 The Motivation

How to get a group of unrelated people to contribute their efforts to solving the tasks set in our two-folded purpose? This section looks at the reasons people might have to contribute their computational powers to a system with a purpose like wePorter. Looking at the way people engage with online video content on platforms like youtube, we identify patterns in their behaviour that can be matched to a task in a human computation system. This behaviour that is characterised by a more active role in multimedia consumption, can be seen as a larger trend in the development of new media. The end of this section indicates how the motivations of users of the wePorter system can link in with this larger trend.

5.4.1 Information Provision through Online Video

Since the proliferation of mobile video recording devices, it has become common practice for large-scale (semi-)public events to be covered in UGVC that gets uploaded to the web. While some are critical[12] to the often heralded democratisation and empowerment of people by the new media production and distribution tools, it is clear that the UGVC at places like youtube attracts a lot of traffic from people looking to be informed about recent events. After all UGVC can have its advantages over traditional media when it comes to video news coverage, especially for unexpected events where traditional media do not have the immediacy of user-generated ‘reports’ recorded by coincidental passersby.

In a recent study as part of the the Pew Research Centers Project for Excellence in Journalism, the most popular video’s from YouTube’s ‘News and Politics’ were analysed for a period of 15 months[19]. The authors of the study exemplify the power UGVC can have in news provision by showcasing frequently viewed videos detailing scenes from the earthquake and subsequent tsunami that hit Japan in March 2011. The week following the disaster, the 20 most viewed news-related videos on YouTube all related to the catastrophic

event and were together viewed more than 96 million times. Most of these videos were recorded by individuals who happened to be in the affected areas when the disaster struck, either uploaded by themselves, or by TV channels who appropriated the content. The study furthermore reports that in the studied period, the most searched term of the month on the YouTube platform as a whole was a news-related event 5 out of 15 months.

While the journalism study above focusses on videos with the ‘News & Politics’ label, information provision about current events might span a larger set of categories. Someone looking for footage in order to get a sense of the atmosphere at a recent music festival or public demonstration, might very well find relevant videos in categories like ‘Entertainment’, ‘Travel & Events’ or ‘Nonprofits & Activism’. Across all of these categories, we are able to find examples of vast collections of UGVC, uploaded in the period following up newsworthy events.

The wePorter system focusses on these kinds of topically related sets that people are currently exploring interactively by browsing from one video to the next. This way of navigation is an intermediary between the goals of *goal-oriented browse* and *unarticulated want*. The apparently aimless browsing is now encapsulated by the event but users still roam freely within this topicalized set of content. By navigating from video to video, watching some and skipping others, users leave attentional traces that give valuable insight into a user’s intentional standpoint.

It is this kind of interactions that are already taking place at a large scale that we like to make use of in the wePorter system. Motivated by the wish to explore informative content, users will instinctively and implicitly contribute their human knowledge to a system that is set up appropriately. This kind of motivation fits the category of ‘implicit work’ as it involves activities that people already engage in for their own reasons[17]. Considering users’ wish to be informed and the interactive way in which they navigate, there is most likely also a factor of enjoyment involved though. We expect though that the more specific motivation of information provision might show to become a valid categorisation for the motivation of people in a HCS as it is a common activity on the web and inherently linked with the hard problem of meaningful interpretation of content.

5.5 The Task

In this section we take a look at how the larger goal of finding intervals of regional interest across time within a single video can be branched out into bite-sized tasks executable by a person in a single interaction. We begin by introducing some conceptual considerations that influenced the interface design. Then a detailed overview of the wePorter web interface is presented. We end with a section focussing on the implementation of the system.

5.5.1 Design Considerations

Below are included several points that have been instructive in the development of the interactive task central to the wePorter system. Some of these point are system requirements, others are more guiding design principles or thoughts that have been inspiring and formative in the development.

wePorter is a web interface

The power of a HCS that relies on data from many interactions is truly unleashed in an online setting, where many people can easily participate and interact. For this goal

alone already, wePorter must be a web-based system. Besides the obvious choice of staying in the realm of the online video content, it makes sense to embed the theoretical explorations of this research in the practicalities of current web technologies. With the ongoing development of technology like HTML5, many new possibilities for a user's web browser are unleashed. The implementation of a research tool concerning online video is a good opportunity for the exploration of the technological possibilities of present day web technologies.

Hyperlinked Multimedia

The power of digital content on the internet lies for a great part in its capacity to be hyperlinked. This alleviates the burden of having to host or recompile pieces of media. Instead, files can simply be played and remixed by reference, leaving their respective sources intact and where they are. In the presentation of their digital video repurposing system 'Diver', Pea et al. indicate the advantages of using a virtual camera controlled by XML-based files that reference parts of source video instead of rendering new video clips[16]:

- "Virtual video clips eliminate the generation of redundant video files, greatly reducing disk storage requirements."
- "No rendering time means vastly improved performance. Users can instantly create and play back dynamic path videos without long video-rendering delays."

An implementation of a system where users interact with content that is dynamically reconfigured in real-time will benefit considerably from a hyperlinked functionality, especially when this takes place in an online setting where bandwidth will be limited.

Localized Interest

In order to elicit users' interest in particular parts within a video, we divide each videos in our initial set of topic-related content into smaller parts when presenting them in a user interaction. Slicing up source videos virtually and playing their parts by reference is made possible by a hyperlinked implementation of the video player.

Interest Elicitation

The user interaction design should enable means to learn about a user's interest in a video at a particular moment in time. This to the purpose of discovering localised regions of interest within single videos.

Considering measures that could indicate how people's interest varies across different videos, an idea that quickly surfaced is that interest might be closely linked to attention. When a piece of content contains something that is interesting to many people, this will most likely result in an increase in views, provided the content is accessible to a variety of people. This simple notion is the idea behind global recommendations that show most popular or 'trending' content. Whether the trending item is a video on a sharing platform or a phrase on a microblogging service, when there is a large number of people attending to it, this is a reason to suspect the item to be of interest, even for people who haven't engaged with it yet.

An obvious limitation of these global recommendations is the lack of personalisation. Personalised recommendations are offered because the content that is globally popular

may not be related to the topics of my interest. For the purpose wePorter is serving however, focus around a particular topic is already in place and are in the first instance mostly interested at picking out the parts that share a high level of interest globally.

Recurrent Interaction

Because of the reliance on data, user's should be able (and encouraged) to engage in the interaction more than once.

Users Between Consumers and Producers

Studies reflecting on new media technology and its incorporation in our everyday life are in recent years often speaking of a media convergence, where multimedia content flows dynamically across multiple media platforms and media audiences take an active, participatory role in their search for entertainment experiences. In his book 'Convergence Culture', Jenkins writes:

"This circulation of media content - across different media systems, competing media economies, and national borders - depends heavily on consumers' active participation. I will argue here against the idea that convergence should be understood primarily as a technological process bringing together multiple media function within the same devices. Instead, convergence represents a cultural shift as consumers are encouraged to seek out new information and make connections among dispersed media content. [...] The term *participatory culture* contrasts with older notions of passive media spectatorship. Rather than talking about media producers and consumers as occupying separate roles, we might now see them as participants who interact with each other according to a new set of rules that none of us fully understands."

Surveying the diverse body of research into interactive TV, Cesar and Chorianopoulos propose a new view that considers content editing, content sharing and content control as an alternative to the more hierarchical 'produce-deliver-consume' paradigm associated with traditional media[4]. The movement from passive consumers to (inter)active contributors indicates new expectations by users of new media applications. The trend of users' more active engagement in new media technology fits well with the approach of interest defined by users' interaction and our proposal of storytelling as structured recommendation.

The Death of the Author, the Birth of Collective Creation

Originally voiced by Roland Barthes, who was contemplating a way of literary writing without the use of a clear narrator. He titled his essay 'the death of the author' to signify the lack of presence of an author in written work following this style of writing. In our aim of reconfiguring interesting sub-clips in a new arrangement, the issue of authorship surfaces in a new context.

Less restrictive forms of digital content licensing, like Creative Commons (CC), mean that it is now possible for content uploaded by its original creator, to be used under specified conditions in a new piece of work by someone else. This kind of licences has been noted to be an important facilitator of research into Human Computation[14]. They make it possible for works not only to be used and remixed by other individuals, but also to be incorporated in algorithmically constructed reconfigurations of user generated content.

Different video platforms are currently offering less-restrictive CC licensing as an integrated part of their services. YouTube currently offers the option of choosing a most basic attribution licence and reports 4 million videos licensed this way [?]. The video platform Vimeo focusses on letting video and animation producers share and showcase their original work. The platform has internalized the use of CC from 2010[?] and many of their users licence their videos such that they can be remixed by others. Figure 5.1 shows that a large part of the licences on the Vimeo platform allow derivatives to be made[?].

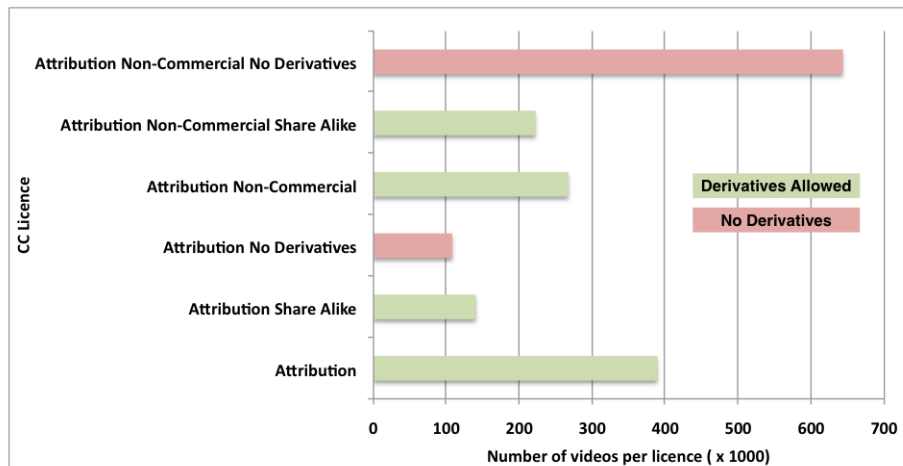


Figure 5.1: Number of videos for each of the Creative Commons licences on Vimeo

5.5.2 The Interface

This section describes the user interface that directs participation of wePorter users towards solving the purpose of distinguishing local intervals of interest within videos. After a conceptual overview of the functioning we include a walkthrough to explain precisely how the interaction takes place.

We hypothesise that interesting parts of content will attract a relatively large amount of attention compared to less interesting parts.

To force a user to make an explicit choice between parts of content for which we would like to elicit their interest, we present two pieces of video playing concurrently and force users to attend to one or the other. In order to get an idea of the variation of interest across a video, ‘source videos’ are divided into smaller ‘video parts’, each of which is presented separately in interactions over time. During this ‘parallel play’ of two video parts we capture the amount of time attended to each of the parts and store this for later analysis.

A Walkthrough

When a user opens the wePorter web interface he is welcomed by a short introduction to the project and successively guided to further instructions explaining the experiment. Introduction header and instructions are shown in ??

In its current version the system uses aggregate counts for the time a user focussed on a particular sub-clip

Parallel player

wePorter

Welcome[About](#)[Get in touch](#)

Welcome to **wePorter**, a research project on [Human Computation](#) in online video storytelling. This website lets you try a new way of watching videos online and at the same time contribute to a research project!

Two videos are presented at once, but you can only 'focus' on one, making the video audible and clearly visible. The two videos sequences are automatically mixed together from parts of different videos. While you are watching, you're submitting data that helps to improve the mix! Please read the full instructions below before you begin.

This website currently only runs in recent versions of Chrome (Sorry no Firefox, Safari or IE). You can download the latest version [here](#).

Besides the main experiment, please also try [two smaller experiments](#) if you haven't already.

Figure 5.2: wePorter Website Header

Instructions✕

In this first experiment, you will be presented with two videos in parallel. This interaction takes 60 seconds.

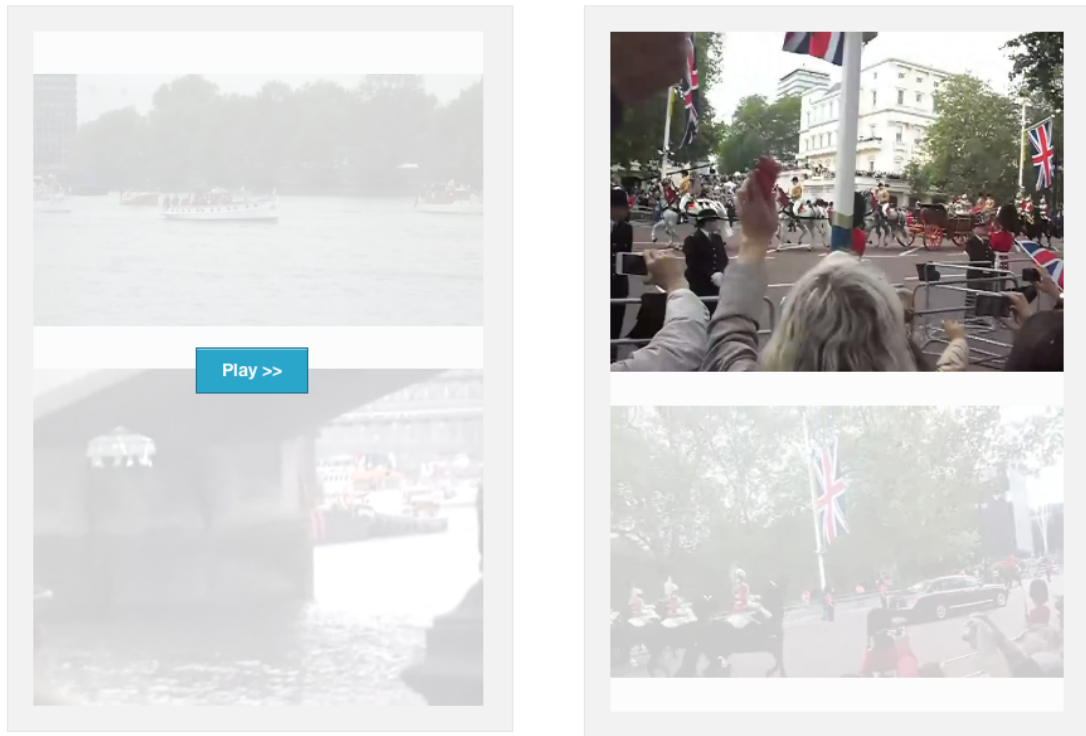
Two videos are presented at once, but you can only 'focus' on one, making the video audible and clearly visible. Focus on a video by moving your mouse over it. The unfocussed video is still dimly visible, allowing you to look what's going on there. You are free to move your mouse over any of the two videos at any time.

The presented videos are recorded at the celebration of the Queen Elizabeth II Diamond Jubilee in May in London.

You can try this experiment as many times as you like.

Once the two videos have loaded, hit **PLAY** to start the videos. (Please reload the page if one of them doesn't load). There will be no option to pause.

Figure 5.3: wePorter Instructions



(a) Upon load

(b) While playing, focus on top video

Figure 5.4: The wePorter parallel play interface

In the web-based interface two video players are presented at once. When a user initiates an interaction session by hitting play, the two videos start playing. Both videos are now visible to the user, but only one of the audio tracks is played. Users control which player’s audio track is being played by moving the mouse cursor over the player they wish to hear. Apart from triggering the audio of a video this way, mouse-over movement also brings focus to a video by displaying it clear, while video out of focus is silent and slightly dimmed.

For every 10th of a second, we record which video a user attends to. Note that we never explicitly ask anyone to point at the video that is most interesting. Users are simply instructed as to how the interface works and then left to explore the video as they like. By recording users’ behaviour this way, we achieve a detailed insight as to which of a pair of videos a user has attended to at what time.

Positioning two video’s one on top of the other, might inflict a bias for users in their attentional behaviour. It might be the case videos on the top are systematically more attended to than videos displayed below. We’ve experimented to see whether such positioning bias effects occur and report on this in section ?? [*TODO ref].

Sequences of video parts

The video sequences played by the two players consist of distinct parts of videos from a topic-centred video database. This means all content that is aggregated in the interaction is focussed around a single topic, in this example the event of the Queen of England’s Diamond Jubilee celebrations in London. The sequences have some other important characteristics: Sequences are formed out of the same number of video parts (in this example 6) Video parts in sequences are all of the same duration (in this example 10 seconds) It follows that all sequences have equal duration (in this example 1 minute).

5.5.3 Forced Feedback

By making data acquisition from interaction an intrinsic part of a human computation system [*need umbrella term], submitting data to the system will become second nature to users and can even take place without the users being aware of it. -i curation through interaction

5.6 Analysis of the Weporter System

5.6.1 Landscapes of Interest

The distinction to be made, between interesting intervals on one hand and less striking parts of a video on the other is not likely to be a very strict one. Afterall, an unedited video captures a single stretch of space and time, so any event that is of particular interest will unlikely have hard cut-off points in time. Rather, if we see interest as a function of time in a particular video, we would expect a somewhat continuous flowing line with spikes every now and then when an interesting event occurs.

Looking at interest at a more global level, aggregating over a large group of users, would perhaps even a more smooth landscape of interest. This kind of data could reveal mountains and valleys that can be used for interest-based segmentation. From the thus segmented parts, the ones with high interest score can be returned as the salient parts within a video.

5.6.2 The death of the author

ref: Barthes - Image music Sound here: a lack of clear narrator of the written story

in wePorter: the lack of a centralised creator, authoring a work consciously and deliberately.

This new form of creation could at first sight be seen as non-authorship, but is rather a collective authorship.

Several recent developments now make possible this collective form of authorship: - the proliferation of tools that enable individual authorship of multimedia content - the increased connectivity of these devices, which gives them the capacity to make the created media accessible to others - the platforms for hosting multimedia (MM) content. This point is strongly related to the previous one. In their development they form a chicken and egg relation. Initially it's hard to imagine one without the other, but once matured, they [influence each other positively]. - the Methods to link, mix, aggregate and modify content online.

5.7 Parallel Play

A new method for preference elicitation in time-based multimedia content.

and a task is set up such that their interactions are recorded as contributions to solving small instances of a larger problem we have two of the three ingredients of our human computation system in place.

5.8 Implementation

argmin
 x

Algorithm 1 My algorithm

```
1: procedure EUCLID( $a, b$ )                                ▷ The g.c.d. of a and b
2:    $r \leftarrow a \bmod b$ 
3:   while  $r \neq 0$  do                                    ▷ We have the answer if r is 0
4:      $a \leftarrow b$ 
5:      $b \leftarrow r$ 
6:      $r \leftarrow a \bmod b$ 
7:   end while
8:   return  $b$                                              ▷ The gcd is b
9: end procedure
```

Algorithm 2 Generate Sequences Random Shuffled

```
1: procedure FILL SEQUENCES                                ▷ The g.c.d. of a and b
2:    $r \leftarrow a \bmod b$ 
3:   while  $r \neq 0$  do                                    ▷ We have the answer if r is 0
4:      $a \leftarrow b$ 
5:      $b \leftarrow r$ 
6:      $r \leftarrow a \bmod b$ 
7:   end while
8:   return  $b$                                              ▷ The gcd is b
9: end procedure
```

Chapter 6

Evaluation

This is the Evaluation section.

Chapter 7

Discussion

Chapter 8

Future Directions

This is the Future Directions chapter.

Creations thus constructed will embody an interesting aspect of collective creation through the merging of content creators, people contributing their human computational capacity and algorithmic aggregation of these components. It will be interesting to see how these developments take form and how they might shape future ideas of authorship.

Chapter 9

Conclusions

This is the conclusions chapter.

Bibliography

- [1] Charles Babbage. On the economy of machinery and manufactures. 1832.
- [2] E. Bruno and D. Pellerin. Video structuring, indexing and retrieval based on global motion wavelet coefficients. *Pattern Recognition, 2002. Proceedings. 16th International Conference on*, 3:287–290 vol. 3, 2002.
- [3] Cathy Casserly. Heres your invite to reuse and remix the 4 million creative commons-licensed videos on youtube. <http://youtube-global.blogspot.co.uk/2012/07/heres-your-invite-to-reuse-and-remix-4.html>, June 2012. Accessed: 15/09/2012.
- [4] Pablo Cesar and Konstantinos Chorianopoulos. The Evolution of TV Systems, Content, and Users Toward Interactivity. *Foundations and Trends® in Human-Computer Interaction*, 2(4):373–95, 2009.
- [5] M. Christel and N. Moraveji. Finding the right shots: assessing usability and performance of a digital video library interface. *Proceedings of the 12th annual ACM international conference on Multimedia*, pages 732–739, 2004.
- [6] M.G. Christel and R.M. Conescu. Addressing the challenge of visual information access from digital image and video libraries. *Proceedings of the 5th ACM/IEEE-CS joint conference on Digital libraries*, pages 69–78, 2005.
- [7] J. Davidson, B. Liebald, J. Liu, P. Nandy, T. Van Vleet, U. Gargi, S. Gupta, Y. He, M. Lambert, and B. Livingston. The YouTube video recommendation system. *Proceedings of the fourth ACM conference on Recommender systems*, pages 293–296, 2010.
- [8] O. De Rooij, C G M Snoek, and M Worring. Query on demand video browsing. *Proceedings of the 15th international conference on Multimedia*, pages 811–814, 2007.
- [9] O. De Rooij, C G M Snoek, and M Worring. Balancing thread based navigation for targeted video search. *Proceedings of the 2008 international conference on Content-based image and video retrieval*, pages 485–494, 2008.
- [10] D.A. Grier. *When Computers Were Human*. Princeton University Press, 2007.
- [11] L Hollink, G P Nguyen, D C Koelma, A Th Schreiber, and M Worring. Assessing user behaviour in news video retrieval. *IEE Proceedings - Vision, Image, and Signal Processing*, 152(6):911, 2005.
- [12] Anna Maria Jönsson and Henrik Örnebring. USER-GENERATED CONTENT AND THE NEWS. *Journalism Practice*, 5(2):127–144, April 2011.

- [13] H. Kuwano, Y. Taniguchi, H. Arai, M. Mori, S. Kurakake, and H. Kojima. Telop-on-demand: Video structuring and retrieval based on text recognition. *Multimedia and Expo, 2000. ICME 2000. 2000 IEEE International Conference on*, 2:759–762 vol. 2, 2000.
- [14] E. Law and L. Von Ahn. Input-agreement: a new mechanism for collecting data using human computation games. *Proceedings of the 27th international conference on Human factors in computing systems*, pages 1197–1206, 2009.
- [15] T. Mei, B. Yang, X.S. Hua, L. Yang, S.Q. Yang, and S. Li. VideoReach: an online video recommendation system. *Proceedings of the 30th annual international ACM SIGIR conference on Research and development in information retrieval*, pages 767–768, 2007.
- [16] R. Pea, M. Mills, J. Rosen, K. Dauber, W. Effelsberg, and E. Hoffert. The diver project: Interactive digital video repurposing. *Multimedia, IEEE*, 11(1):54–61, 2004.
- [17] A.J. Quinn and B.B. Bederson. Human computation: a survey and taxonomy of a growing field. *Proceedings of the 2011 annual conference on Human factors in computing systems*, pages 1403–1412, 2011.
- [18] Pablo Cesar Dick C A Bulterman Vilmos Zsombori Ian Kegel Rodrigo Laiola Guimaraes. Creating Personalized Memories from Social Events: Community-Based Support for Multi-Camera Recordings of School Concerts. Technical report, August 2011.
- [19] Tom Rosenstiel and Amy Mitchell. YouTube & the News. Technical report, July 2012.
- [20] Y. Rui, T.S. Huang, M. Ortega, and S. Mehrotra. Relevance feedback: A power tool for interactive content-based image retrieval. *Circuits and Systems for Video Technology, IEEE Transactions on*, 8(5):644–655, 1998.
- [21] A.F. Smeaton, P. Over, and W. Kraaij. Evaluation campaigns and TRECVID. *Proceedings of the 8th ACM international workshop on Multimedia information retrieval*, pages 321–330, 2006.
- [22] AM Turing. Computing machinery and intelligence. *Mind*, 1950.
- [23] Marian F Ursu, Vilmos Zsombori, John Wyver, Lucie Conrad, Ian Kegel, and Doug Williams. Interactive documentaries. *Computers in Entertainment*, 7(3):1, September 2009.
- [24] S. Xu, H. Jiang, and F. Lau. Personalized online document, image and video recommendation via commodity eye-tracking. *Proceedings of the 2008 ACM conference on Recommender systems*, pages 83–90, 2008.
- [25] B. Yang, T. Mei, X.S. Hua, L. Yang, S.Q. Yang, and M. Li. Online video recommendation based on multimodal fusion and relevance feedback. *Proceedings of the 6th ACM international conference on Image and video retrieval*, pages 73–80, 2007.
- [26] M. Yang, B.M. Wildemuth, and G. Marchionini. The relative effectiveness of concept-based versus content-based video retrieval. *Proceedings of the 12th annual ACM international conference on Multimedia*, pages 368–371, 2004.

- [27] R. Zhou, S. Khemmarat, and L. Gao. The impact of YouTube recommendation system on video views. *Proceedings of the 10th annual conference on Internet measurement*, pages 404–410, 2010.
- [28] V. Zsombori, M. Frantzis, R.L. Guimaraes, M.F. Ursu, P. Cesar, I. Kegel, R. Craigie, and D. Bulterman. Automatic generation of video narratives from shared UGC. *Proceedings of the ACM Conference on Hypertext and Hypermedia*, pages 325–334, 2011.