

<b>Noname manuscript No.</b> (will be inserted by the editor)
--

---

## Crowdsourcing Authoring of Sensory Effects on Videos

Marcello Novaes de Amorim · Estêvão Bissoli Saleme · Fábio Ribeiro de Assis Neto · Celso A. S. Santos · Gheorghita Ghinea

Received: date / Accepted: date

**Abstract** Human perception is inherently multi-sensorial involving five traditional senses: sight, hearing, touch, taste, and smell. In contrast to traditional multimedia, based on audio and visual stimuli, mulsemmedia seek to stimulate all the human senses. One way to produce multi-sensorial content is authoring videos with sensory effects. These effects are represented as metadata attached to the video content, which are processed and rendered through physical devices into the user's environment. However, creating sensory effects metadata is not a trivial activity because authors have to identify carefully different details in a scene such as the exact point where each effect starts, finishes, and also its presentation features such as intensity, direction, etc. It is a subjective task that requires accurate human perception and time. In this article, we aim at finding out whether a crowdsourcing approach is suitable for authoring coherent sensory effects associated with video content. Our belief is that the combination of a collective common sense to indicate time intervals of sensory effects with an expert fine-tuning is a viable way to generate sensory effects from the point of view of users. To carry out the experiment, we selected three videos from a public mulsemmedia dataset, sent them to the crowd through a cascading microtask approach. The results showed that the crowd can indicate intervals in which users agree that there should be insertions of sensory effects, revealing a way of sharing authoring between the author and the crowd.

---

M. N. Amorim · E. B. Saleme · F. R. A. Neto · C. A. S. Santos  
Federal University of Espírito Santo, Av. Fernando Ferrari, 514, 29075-910 - Vitória-ES, Brazil  
Tel.: +55 27 4009-2324, Fax: +55 27 4009-2850  
E-mail: [cellonovaes@gmail.com](mailto:cellonovaes@gmail.com), [estevaobissoli@gmail.com](mailto:estevaobissoli@gmail.com), [fabio.ribeiro.neto@gmail.com](mailto:fabio.ribeiro.neto@gmail.com), [saibel@inf.ufes.br](mailto:saibel@inf.ufes.br)

G. Ghinea  
Brunel University London, Wilfred Brown Building 215, Kingston Lane, Uxbridge, Middlesex UB8 3PH  
Tel.: +44 (0)1895 274000, Fax: +44 (0)1895 232806  
E-mail: [george.ghinea@brunel.ac.uk](mailto:george.ghinea@brunel.ac.uk)

## 1 Introduction

Mulsemmedia has been designated as AV (AudioVisual) content increased with other non-traditional sensory objects such as olfactory, gustatory and haptic [16]. Throughout the last decade, researchers have been exploring this coherent combination of human senses to enhance the Quality of Experience (QoE) of users in mulsemmedia applications [38, 29, 36, 40, 13, 26, 1]. Nonetheless, mulsemmedia applications face a wide spectrum of research challenges, many of which are by now traditional in the multimedia community. Of these, we mention rendering, distribution, adaptation, sensory cue integration, building mulsemmedia databases, usability, compliance, as well as a lack of rapid prototyping tools [9, 16, 32].

Indeed, many challenges stem from non functional requirements and, in this context, inter-operability is primordial. Towards this end, the MPEG-V standard emerged to provide an architecture and specifications for representing sensory effects [20]. An AV content annotated with MPEG-V Sensory Effects Metadata (SEM) should be able to be reproduced efficiently in many different mulsemmedia systems even with actuators from different brands. The process of producing interoperable mulsemmedia metadata involves making an MPEG-V compatible XML file that contains entities describing different sensory effect presentation features (beginning and the end of each effect time span, its intensity, fading and so forth). This can be done with the help of an authoring tool, which enhances considerably the efficiency of the process overall [37, 21, 7, 33]. This tool allows authors to abstract the difficulty of writing an XML application, through an intuitive graphical interface whereby they can pick up a movie scene, see what they would feel whilst immersed in the scene, and then, arrange the sensory effects on it. Furthermore, researchers have started developing tools and methods to automatically generate MPEG-V SEM from video content [22, 28]. However, there is much research to be done to produce guidelines for authoring, editing and creating mulsemmedia applications. Moreover, the difficulty of capturing many different aspects and turning it into trustful sensory effect metadata remains an issue currently. Both are challenges which we address in this paper. Accordingly, in this article, we explore whether a crowdsourcing approach can generate interoperable metadata and boost the process of authoring mulsemmedia content.

Despite not being complex in terms of manipulating a tool, the authoring of AV content with sensory effects requires some skills from the author, who should be able to identify and capture different details in scenes, such as the exact point when an effect starts and finishes, the type of effect to be presented and other details. This is a manual and subjective process usually requiring accurate human perception, time and the ability to produce a coherent combination of sensory effects and AV content.

In the process of authoring mulsemmedia content, individuals are subject to natural bias due to their unique prior experiences and the QoE of users is subject to what they are feeling. Thus, we believe that the combination of a

collective common sense to indicate time intervals of sensory effects with an expert fine-tuning is a viable alternative to efficiently generate metadata.

Assuming that the purpose of the experiment is use the crowd to authoring of sensory effects, the expected result is the insertion of *Wind* effects and *Vibration* in the intervals in which the crowd agrees that it makes sense to insert them.

Our approach is based on the Galton concepts for Wisdom of Crowds [15]. Following the same principle, Chowdhury et al. [8] built the Know2Look system and obtained satisfactory results by using common sense to determine valid annotations in media content.

Accordingly, in the study reported here, we selected three videos from a public mulsemmedia dataset [39]. Next, we sent them to workers, recruited on the MicroWorkers<sup>1</sup> platform, through a cascading microtask approach. Then, we analyzed the gathered metadata from the crowd. Finally, we placed the results in parallel to the dataset to make a discussion about our approach and this theme. Early results revealed that the sensory effects generated by the crowd can be used to validate author insertions as well as to supplement authorship with new sensory effects that members of the public agree should be added.

The remainder of the text is organized as follows. Section 2 brings other works related to the authoring of mulsemmedia content and crowdsourcing approaches for multimedia authoring and annotation. Section 3 presents the workflow that guides our crowdsourcing approach for authoring mulsemmedia content. Section 4 describes the tools that support our approach. Section 5 presents the evaluation of the work. Section 6 discusses the results. Finally, Section 7 concludes the article and leads to future works.

## 1.1 Scope of this Work

Initially it is important to make clear that this article aims to present a crowdsourcing approach to sensory effects video authoring. Strategies or new features for its annotation are out of the scope. In our approach, workers make trivial annotations on videos, so the tools used to collect them are simple HTML documents that contain players and additional controls to indicate instants and intervals. Subsequently, the collected annotations are filtered, grouped, and aggregated to produce results that represent the common sense from the crowd. There are models and definitions that aim to standardize the production of annotated media, such as the canonical process presented by Hard et al. [18] for semantically annotated media production. However, the model presented in this work are related to the cascade microtasking process of crowdsourcing, not to the annotation process itself. The annotation collected from the crowd is straightforward and used as an input to the aggregation method, which in turn, process the contributions and generates the final results.

---

<sup>1</sup> MicroWorkers platform available at <http://ttv.microworkers.com>

Complementarily, it is not within the scope of this article to propose a presentation model for mulsemmedia videos, similar to that presented by Sadallah [30] for hypervideo on the web. Instead, the outcome will be exported to interoperable and sharable formats, such as MPEG-V, which allows results accessed from compatible systems.

## 2 Related work

Mulsemmedia authoring tools have been developed for almost a decade. SEVino [37], SMURF [21], RoSE Studio [7], and Real 4D studio [33] are examples of tools that support authors in adding sensory effects, usually represented as MPEG-V metadata, to AV contents. Players compatible with MPEG-V such as Sensorama [6], SEMP [37], Sensible Media Simulator [21], Multimedia Multisensorial 4D Platform [4], and PlaySEM [31] are able to reproduce this kind of authored content. Thus, all of these tools shape an ecosystem for delivering and rendering mulsemmedia content.

The work of Kim et al. [22], and more recently the work of Oh and Huh [28], represent an attempt to automatically generate interoperable mulsemmedia metadata. The authors argue that this method can speed up the authoring process, helping the industrial deployment of mulsemmedia services. Kim et al. [22] proposed a method and an authoring tool to extract temperature effect information using the color temperatures of video scenes and generate MPEG-V SEM. The authors found that the automatically generated temperature effects enhanced the level of satisfaction of the users through a subjective evaluation. However, its limitation to generate only one effect relies on the recurrence of manual tools to add other effects. Oh and Huh [28] proposed a similar approach to automatically generating MPEG-V SEM based on the motion of an object included in a media. Akin to the approach of Kim et al. [22], it is limited to the temperature effect information, automatically extracted from the color temperatures of video scenes. Furthermore, the authors did not show the results of the method, which makes it difficult to evaluate its efficiency.

Interoperability is an important issue to be addressed in systems for authoring and annotating media. This problem is often addressed in works in this area, such as Sadallah et al. [30], which presents a high-level component-based model for specifying annotation-driven systems focused on hypervideo. In this context, it is interesting to adopt models and standards that promote the interoperability of generated metadata, such as the canonical process proposed by Hardman et al. [18] for the production of semantically annotated media. The metadata interoperability allows its use in different applications, including by automatic methods, as can be observed in Ballan et al [3] that has surveyed works focused on the detection and recognition of events for semantic annotation of videos.

Over the past years, crowdsourcing has been applied to many studies involving multimedia content processing, such as for generating video annotations [24,11], image descriptions [14,34], real-time captioning [23], text an-

notations [12,42], and audio annotations [25,35]. Recently, See and Chat [5] demonstrated how to automatically generate annotations in user comments on images published on Flickr<sup>2</sup>.

With regard to the generation of complex video metadata using crowdsourcing methods, the work of Cross et al. [10] presented the VidWiki system, which is a complex application to improve video lessons. This system requires that workers edit video scenes by adding annotations, including complex annotations such as LaTeX equations. It also requires that the recruited workers have previous knowledge about LaTeX. Another crowdsourcing complex video annotation system was proposed by Gottlieb et al. [17], which achieved geo-location annotations for random videos from the web by requiring the workers to perform searches on the internet, use encyclopedias and provide annotations in the specific format for GPS coordinates.

In relation to the crowdsourcing processes for multimedia authoring, it is possible to cite the work of Amorim et al. [2], that used a process composed by cascade microtasks to generate interactive multimedia presentations from plain videos. In this work the audience was responsible for identifying the points of interest to be enriched, as well as making available, selecting and positioning the media content in the video to generate the multimedia composition.

As far as authoring of mulsemmedia content based on a crowdsourcing approach, we did not find related studies dealing with it employing a multi-stage crowdsourcing process, in which activity diagrams represent the process workflow and in which quality will be managed during the different execution stages of this process.

### 3 A Crowdsourcing Approach for Authoring of Sensory Effects

Our process for mulsemmedia content authoring is composed by a sequence of microtasks, so as to be possible to use the workforce of a plethora of unskilled workers in terms of mulsemmedia. In fact, this process can be viewed as a generic solution that is tailored to different types of crowdsourcing annotation projects. Figure 1 presents the three phases - *Preparation*, *Execution*, and *Conclusion* - of this generic process.

The *Preparation Phase* deals with the definition of AV content to be annotated, the source of the contributions for each task, the monetary resources to pay workers, and the tools used to them for performing the tasks.

The *Execution Phase* deals with the process workflow, which is generically is described in the form of a complete algorithm that controls task flows within time, cost and quality constraints, to reach a desired end result. In our approach for mulsemmedia authoring, crowdsourcing tasks are performed in cascade, a sequence of several similar stages with each stage processing the output from the previous stage. In addition, all task executions are associated to the same sequence of steps: *Selecting Workers*, *Collecting Contributions*, *Filtering*

---

<sup>2</sup> Flickr available at <https://www.flickr.com>

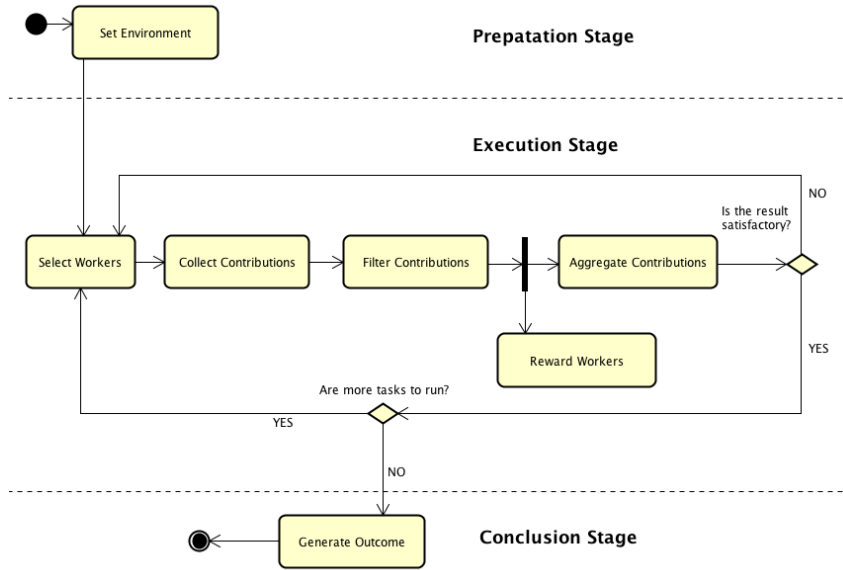


Fig. 1: Generic process template's workflow.

*Contributions, Reward Workers, and Aggregating the Filtered Contributions.* The results generated by the aggregation method for a task provides the input for the next one. Once the output of the aggregation method is satisfactory, the process advances to the next stage in the cascade for further processing. Otherwise, the current task in the workflow must be restarted to select new workers, collect new individual results and update the aggregated task results. *Reward Workers* are located after the *Filter Contributions* to cover cases where the *Owner* has decided to pay only for valid contributions. Thus, if the process requires another task, simply follow that model, defining the annotation tool, the filters and the aggregation method, and insert it into the process.

In the *Conclusion Phase*, the end result is produced and evaluated using an specific method defined in the project.

The generic process described in this section process can be used to create crowdsourcing workflows to coordinate the crowd through a sequence of tasks, managing their dependencies, and bringing together intermediate results produced by the workers as in the case of Figure 2.

In fact, a workflow represents an effortless way to understand the whole process since from hiring workers until processing the results provided by the crowd. In addition, as stated by Assis Neto and Santos [27], crowdsourcing workflows are context-oriented and they should establish not only how the process activities will be performed by the crowd but also how the quality will be managed through the different execution stages of the crowdsourcing process.

### 3.1 The Workflow for Mulsemedia Authoring

The workflow presented in Figure 2 represents our crowdsourcing process for mulsemmedia authoring based on AV content annotation with MPEG-V SEM. This workflow consists of a set of tasks performed by actors performing four roles: *Owner*, *Crowdsourcing Platform (CS Platform)*, *Crowdsourcing Process Manager (CSPM)*, and *Worker*.

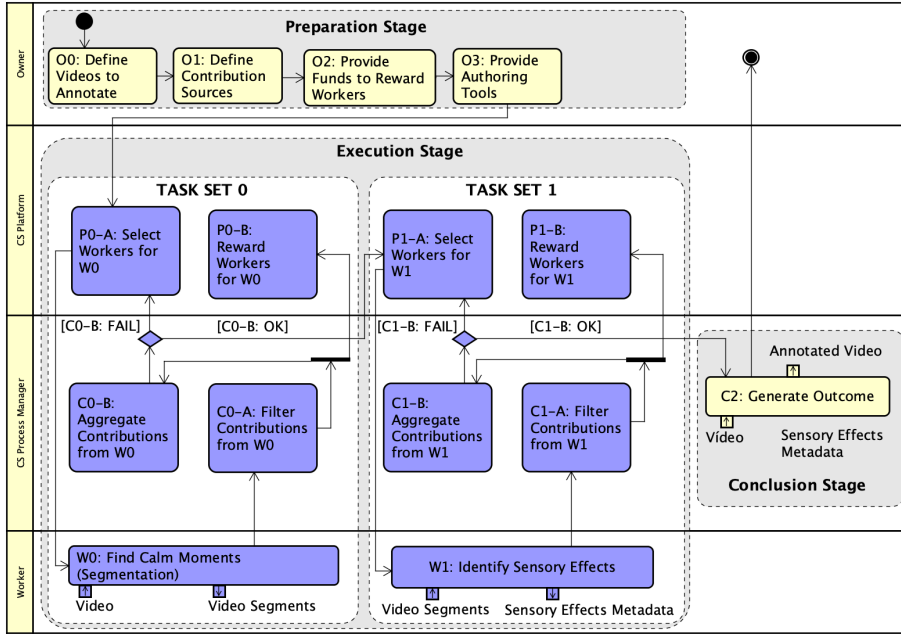


Fig. 2: Crowdsourcing process's workflow for mulsemmedia content authoring.

The process begins with the *Preparation Stage* (see Figure 1) in which the *Owner* sets the environment to start the process. In our view, the *Owner* is someone who works with a crowdsourcing management team, composed of experts in the field, and responsible for specifying the technical requirements as well as qualified personnel to create the tasks.

At the beginning of this stage, the owner performs the activity ***O0: Define Videos to Annotate***, registering the videos to be annotated. Next, he/she must then perform the activity ***O1: Define Contribution Sources***, in which he/she chooses whether to use a commercial crowdsourcing platform or other mean to reach workers to collect their contributions. When choosing to use a commercial crowdsourcing platform, it is necessary to deposit funds to reward workers, this is done in activity ***O2: Provide Funds to Reward Workers***. Completion of the *Preparation Stage* is reached when the activity ***O3: Provide Annotation Tools*** is concluded and the campaigns are created in the

*CS Platform*. In our case, the campaigns correspond to the crowd tasks **W0: Find Calm Moments** and **W1: Identify Sensory Effects**.

The *Crowdsourcing Platform* is the source of contributions, it is responsible for the activities **P0-A and P1-A: Select Workers for W0 and W1**, and **P0-B and P1-B: Monetary Reward Payment for W0 and W1**.

The *Crowdsourcing Process Manager (CSPM)* represents a person, or a team, responsible for monitoring the process and analyzing the state of contributions to decide when a task should be closed, as well as initiating the generation of results for each task, and stop and start tasks.

The *CSPM* is responsible for activities that produce partial results and compile into the outcome.

In the activities, **C0-A and C1-A: Filter Contributions from W0 and W1** reliability filters are applied over the collected annotations from the crowd, so activities **C0-B and C1-B: Aggregate Filtered Contributions from W0 and W1** process reliable contributions to construct the results. Finally, after all the partial results are made, *CSPM* executes activity **C2: Generate Outcome** to export the annotated video to the desired format.

The *Workers* are responsible for providing the information required to produce the outcome. They are responsible for performing the annotation tasks by execute the activities **W0: Find Calm Moments** and **W1: Identify Sensory Effects**.

Responsibilities and activities for each role are summarized in Table 1.

Table 1: Responsibilities and activities for each role.

Role	Responsibilities	Activities
<b>Owner</b>	Setting up the environment; Provide funds to pay Workers.	O0
		O1
		O2
		O3
<b>Crowdsourcing Platform (CS Platform)</b>	Recruiting Workers; Intermediate payments.	P0-A
		P0-B
		P1-A
		P1-B
<b>Crowdsourcing Process Manager (CSPM)</b>	Manage the process workflow; Initiate Filtering and Aggregation; Control transitions between tasks.	C0-A
		C0-B
		C1-A
		C1-B
		C2
<b>Worker</b>	Execute the annotation tasks.	W0
		W1

#### 4 CrowdMuse: Crowdsourcing Multimedia Authoring System

The *CrowdMuse* system has been developed to support our crowdsourcing approach. One of the most important characteristics of the system is its ca-



capacity of distribute tasks to workers from various sources, such as commercial crowdsourcing platforms, internal teams, and social networks.

*CrowdMuse* follows a component-based approach to manage the complexity of a mulsemmedia annotation problem by breaking it down into smaller and physical manageable modules. The modules *Server*, *Client*, and *Persistence* in Figure 3 are the units of implementation of the *CrowdMuse* system and are assigned areas of functional responsibility. In addition, the work interfaces in the system were constructed as simple HTML-5 documents, which simply render information coming from the *Server Module* and send back contributions.

Another advantage of the *CrowdMuse* architecture is that even when using commercial crowdsourcing platforms, the collected data is stored only in the system database. Moreover, this system is responsible for controlling the execution flow of the tasks, managing the items that must be annotated, generating the jobs that must be executed and distributing these jobs among the workers.

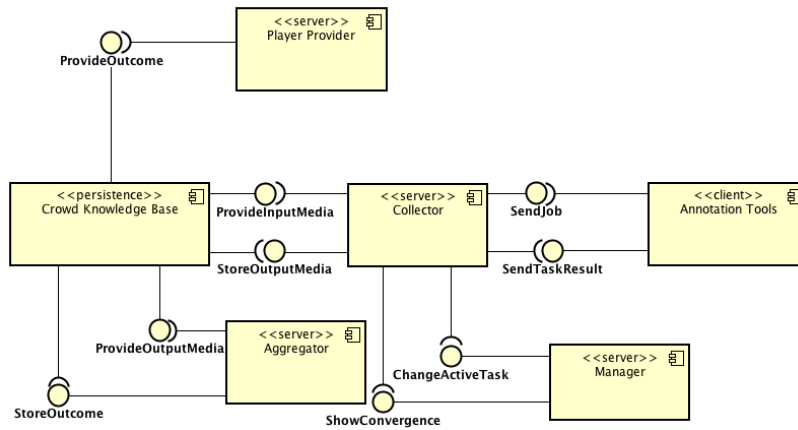


Fig. 3: CrowdMuse system components and communication interfaces.

#### 4.1 The Persistence Module

In the *Persistence Module*, the *Crowd Knowledge Base* component sends to the server module the information required to render the job requests to workers and the content needed to present the result to the users. Likewise, contributions produced by workers were sent directly from the collector to the *Crowd Knowledge Base* component, in which they were stored directly in the database without having to go through the external crowdsourcing environment.

The *Aggregator* component also communicates directly with the persistence. The *Aggregator* retrieves the collected contributions of a task set and,

after the aggregation process, sends the result to be stored in the database, to be used as input to the next task, again maintaining data privacy because it does not need to be stored in an external environment.

## 4.2 The Server Module

The server module is responsible for distributing the jobs, managing contributions, and controlling the active task in order to execute the process workflow. This module is composed of four components: *Manager*, *Collector*, *Aggregator*, and *Player Provider*.

- **Manager** is the module responsible for controlling the enrichment process. It is related to the task-to-task transition tool, as well as the tool to monitor the current state of tasks and trigger aggregation methods. Management functionality is accessible through the management interface.
- **Collector** provides the annotation tool with information about the item to be annotated, therefore, it renders the job's interface used by the worker to perform the task. Also, this component is responsible for gathering the information provided by the worker after the execution of the task and sends them to the persistence.
- **Aggregator** applies reliability filters over the worker's contributions and processes the valid annotations in order to generate the result for each task. The aggregation methods are based on convergence analysis to produce collective results.
- **Player Provider** delivers the process outcome that consists of mulsemmedia annotated videos. This outcome can be exported and visualized in players able to reproduce these effects, such as SEVino [37] tool.

## 4.3 The Client module

The *Client Module* manages the communication interfaces involving workers and other users. This module presents templates that generate visualization data according to a description. For each task of the process, the client must render a specific annotation tool for the worker. Thus, it is possible to keep the server accessible through the *Collector* and *Player* components, and therefore templates can be stored from anywhere. This allows contributions to be collected from different workers at the same time. Moreover, a model is selected and the Persistence module is queried to obtain the necessary data to render the desired interface according to the desired data visualization.

## 4.4 The Public Interfaces

The communication between the modules of the *CrowdMuse* system occurs through the public interfaces represented in Figure 3 and detailed as follows:

- **Change Active Task:** The *Owner* sends a request to the Server to set the currently active task.
- **Show Convergence:** The *Manager* displays to the owner the current convergence state for the active task.
- **Provide Media Input:** To generate each job, the *Collector* receives an entry from the *Persistence* component.
- **Send Job:** The *Collector* sends a job to a worker who sees the task through the *Client* and executes it.
- **Send Task Result:** The *Client* sends to the *Collector* the annotation made by the Worker.
- **Store Media Input:** The *Collector* sends workers contributions to the *Persistence*, that stores it in the *Crowd Knowledge Base*.
- **Provide Output Media:** The *Persistence* send to the *Aggregator* all the contributions collected related to a task.
- **Store Outcome:** The *Aggregator* stores the entries received from the aggregation process. The generated outcome can be provided as input to the next task.
- **Provide Outcome:** The *Persistence* module provides the outcome of the crowdsourcing project (i.e., a set of MPEG-V SEM) to be rendered with the video content.

#### 4.5 Considerations

The CrowdMuse system can be freely used and modified to serve different crowdsourcing applications with a focus on authoring of mulsemmedia and other kinds of multimedia content. In the next section, we will present a case study demonstrating the use of CrowdMuse for crowdsourcing authorship of mulsemmedia content according to our approach.

### 5 A Case Study on Crowdsourcing Authoring of Mulsemmedia Content

A case study concerning the crowdsourcing authoring approach proposed here was carried out using three from a public mulsemmedia dataset [39], referenced in the paper as *Babylon A.D.*, *Formula 1*, and *Earth*. As stated by Timmerer et al. [36], these three videos have obtained the highest MOS in their QoE experiments which were performed over this same dataset and that is the main reason for this choice. Despite not having enough evidence, we assumed that workers could perceive sensory effects easier in videos with high MOS than in random videos, and thereby, give a clearcut contribution.

The reference mulsemmedia dataset contains also information about the intensity of some sensory effects. However, because there was no homogeneity between the audio and video equipment used by the workers, it was decided not to request that they observe the intensity of the effects, only the intervals

at which they should be inserted. In addition, we noticed that in the dataset of Waltl et al. [39] there are annotations of effect with very subtle intensity, and we chose not to consider them for this experiment, believing that a more specialized work would be needed to accurately record the subtle effects.

We decided to collect the same kind of effects associated by of Waltl et al. [39] to the selected videos, that is *Wind* and *Vibration* effects. The metadata associated to lighting information, also annotated in the videos, will be not considered in our experiment once it is set to be auto-extracted according to the brightness and color information of the video frames.

Table 2 presents the main characteristics of the three videos annotated with sensory effects used in our evaluation.

Table 2: AV content annotated with sensory effects used in the experiment.

Video	Resolution (WxH@fps)	Bit-rate (Mbit/s)	Duration (s)	Wind	Vibration
<b>Babylon A.D.</b>	1280x544@24	6.81	118.42	10	8
<b>Earth</b>	1280x720@25	6.90	66.00	7	1
<b>Formula 1</b>	1280x720@25	5.40	116.2	11	4

To conclude this section, we come to the first question posed in the introduction of the paper: Is the crowd is capable of producing coherent and cohesive set of sensory effects related to the AV content processed by each worker individually? Other questions have to do with effort and quality of the content produced in a crowdsourcing process.

## 5.1 Setting the Environment

According to the workflow of Figure 2, the *Owner* should perform four activities to set up the environment before beginning to collect contributions from the crowd. The first activity is ***O1: Set videos to annotate*** and consists of selecting the videos that should be annotated. These videos should be uploaded and set to public. In the activity ***O2: Define Contribution Sources***, the *Owner* chooses if the contributions will be collected from a contracted crowd using a commercial platform or the workers will be volunteers or members of internal groups.

The *Owner* may also need to create a campaign on the crowdsourcing platform (Microworkers, in our case) for each task in the process workflow. To create a campaign, the *Owner* must perform the activity ***O3: Provide Funds to Reward Workers*** to ensure the means to pay workers, and the activity ***O4: Provide Annotation Tools*** in which the annotation tool that will be used to perform the task is sent to the crowdsourcing platform.

## 5.2 Crowd Definition

As mentioned before, *Microworkers* performs the role of *CS Platform* in the project workflow of Figure 2. Thus, this crowdsourcing platform is responsible for recruiting and paying the workers, although *CrowdMuse* is compatible with other commercial platforms, such as Amazon’s Mechanical Turk (AMT)<sup>3</sup> and CrowdFlower<sup>4</sup>.

*Microworkers* proposes different models for a crowdsourcing project. Initially, one can choose between starting a basic campaign or using contracted groups. In a basic campaign, all registered workers in the platform can see the task in the job menu and work on it. A campaign that uses hired groups allows the *Owner* to select the crowd by choosing groups of workers with the desired profile. In addition, it is possible to create lists with workers who have made good contributions before, so they can be recruited to work on other tasks.

One of the characteristics of our approach is that it use very simple tasks and unskilled workers. The tasks were launched as campaigns that used contracted groups, to increase the chance of the workers who contributed to a task also participating in others.

A group of moderate size was chosen so the contributions were made quickly. The group chosen is identified as *Data Services* in *Microworkers* platform, with 1285 potential workers to accept the jobs. Some groups were composed of workers who only accepted tasks that offered slightly larger payment, but considering the chosen group, it was feasible to offer a payment of 0.05 USD per task.

## 5.3 Method

As shown in the workflow of Figure 2 the crowdsourcing process for authoring mulsemmedia content is based on two microtasks executed in cascade. Each microtask is executed as a complete task-set construction composed of three sequential main activities: (i) contributions collection, (ii) filtering, and (iii) aggregation. Hence, the individual contributions are collected from each worker through a specific tool required for executing the assigned task. In the sequence, the contributions are filtered and clustered using the aggregation function to extract the results of the microtask execution.

The next two subsections will describe the two tasks (Find Calm Moments and Identify Sensory Effects) that composed the crowdsourcing process for mulsemmedia authoring.

---

<sup>3</sup> AMT - <https://www.mturk.com>

<sup>4</sup> CrowdFlower - <https://crowdflower.com>

#### 5.4 The First Crowd Task: Find Calm Moments (Segmentation)

The objective of the first task was to segment the video in such a way that the sensory effects were not fragmented by more than one segment, that is, the effects contained in a segment should be completely contained in it. In this task, one of the three selected videos was displayed to the worker who should indicate the instants that he/she thought there is no *Wind* or *Vibration*, pretending that he/she was immersed in the environment of the movie.

##### 5.4.1 Contributions Collection

To support the first microtask the tool shown in Figure 4 has been created. In this tool, the video to be processed is on the bottom of the window and the buttons used by the worker to determine the calm moments on the video as well as the task instructions are on the top. As discussed, the instants pointed out by the crowd are used later for content segmentation.

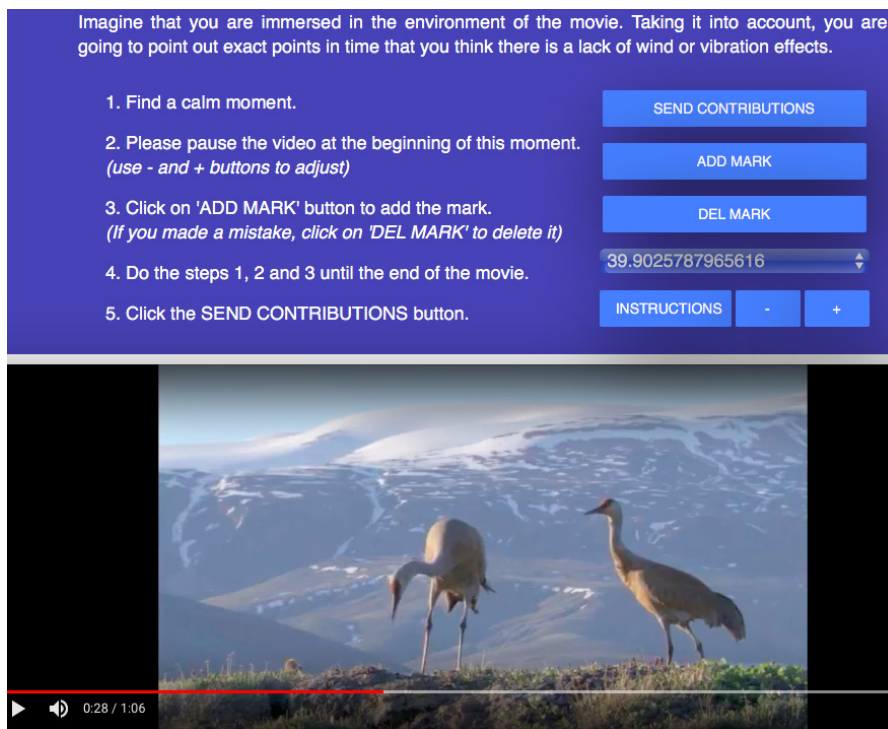


Fig. 4: Tool for identifying calm moments in a video.

In each contribution, a worker could supply as many time marks as he/she wanted, each mark representing the initial instant of a calm moment in the

video. In this task execution, 23 contributions were obtained that provided 113 time-points for the *Babylon A.D.* video, 17 contributions that provided 213 points for the *Formula 1* video, and 21 contributions that provided 108 points for the *Earth* video.

#### 5.4.2 Filtering Contributions

The collected contributions were filtered according to (i) the number of time-points provided and (ii) the proximity of these time-points.

With each task, workers could annotate multiple marks in the video segment. Thus, the first quality criterion was to calculate the average number of marks received by each segment and to discard the contributions containing differing amounts of marks. Contributions with less than 50% or more than 200% of the average number of marks were discarded. In addition, very close marks of the same worker contribution were discarded. It was established that the segments should be at least 0.5 second in length. Thus, when a worker provided two separate marks for less than 0.5 second, a first annotation was ignored for assuming that an update occurred and the worker forgot to delete the previous one.

Summarizing, a total of 68 over 113 time marks remained after the filtration stage for *Babylon A.D.* video. For *Earth* video, 59 over 108 time marks and for *Formula 1* video, 179 over 213 time marks were delivered to the aggregation stage.

#### 5.4.3 Aggregation

The aggregation process for the first microtask is based on the grouping of the contributions so that each group contains suggestions from the crowd regarding the same calm instant in AV content. Because marks are point values that represent instants of the video, the algorithm used is based on neighborhood grouping. This strategy assumes that the distance between two marks referring to the same calm moment tend to be closer than the marks for consecutive calm moments.

In our aggregation strategy for this task, the marks were sorted in ascending order, and a  $\Delta$  value was calculated that represents the average distance between the consecutive marks provided. This  $\Delta$  was then used as the threshold for the grouping. When the distance between one time-point contribution and the next one is greater than  $\Delta$ , a new group is started.

At the end of this stage, each group obtained represents the initial instant of a calm moment for which the crowd agreed to exist. Therefore, time-points that do not fit into any group were discarded.

The video segments are determined using the calm moments defined by the crowd. Each video segment to be annotated with sensory effects is associated to the time interval between two calm moments. With this strategy, we aimed to obtain segments that contain *Wind* and *Vibration* effects without being fragmented by more than one video segment. In this way, a total of 11, 12,

and 15 segments were obtained for the *Babylon A.D.*, *Earth*, and *Formula 1* video, respectively (see Table 3).

Finally, the segments to be annotated with sensory effects were determined are used as input for the second microtask in our authoring approach.

Table 3 shows that, out of a total of 61 contributions, 434 suggestions of calm moments were observed by the crowd. After filtering, that produces 306 instants, whilst 38 segments were obtained after running the aggregation stage.

Table 3: Contributions and results for the first task.

Video	Contributions	Calm Moments	Filtered	Segments
Babylon A.D.	23	113	68	11
Earth	21	108	59	12
Formula 1	17	213	179	15
<b>Total</b>	<b>61</b>	<b>434</b>	<b>306</b>	<b>38</b>

### 5.5 The Second Crowd Task: Identify Sensory Effects

The second task asks workers to provide subjective information, aimed at obtaining ranges in which *Wind* or *Vibration* effects should be pointed out. We re-enforced the workers to provide the maximum time ranges they could and be trustful in an attempt to receive more reliable contributions in this task. We did not ask for intensity because we believe it is a fine-tuning task that requires expert skills such as fade-in and fade-out. However, we advised the workers to create a new range of the same effect if they realized a change in intensity. The input of this second task was the set of 38 video segments produced in the previous task, as detailed in Table 3.

The segments resulting from the aggregation of the contributions collected in the first microtask were delimited by moments of total absence of effects, so that it is possible that in these second microtasks more than one insertion of effect within each segment is identified. This occurs in cases where two inserts of high intensity effects are separated by a low intensity sensory effect, without ceasing completely the effect. In this way, it is possible and acceptable that the number of sensory effects identified in the videos at the end of this task is greater than the number of segments received as input.

#### 5.5.1 Contributions Collection

The tool depicted in the Figure 5 has been implemented to support the collection activity related to the second task. The look-and-feel is very similar to the tool presented in Figure 4. The instructions followed by the workers and buttons for correctly performing the task are presented, at the top, and the analyzed video, at the bottom.



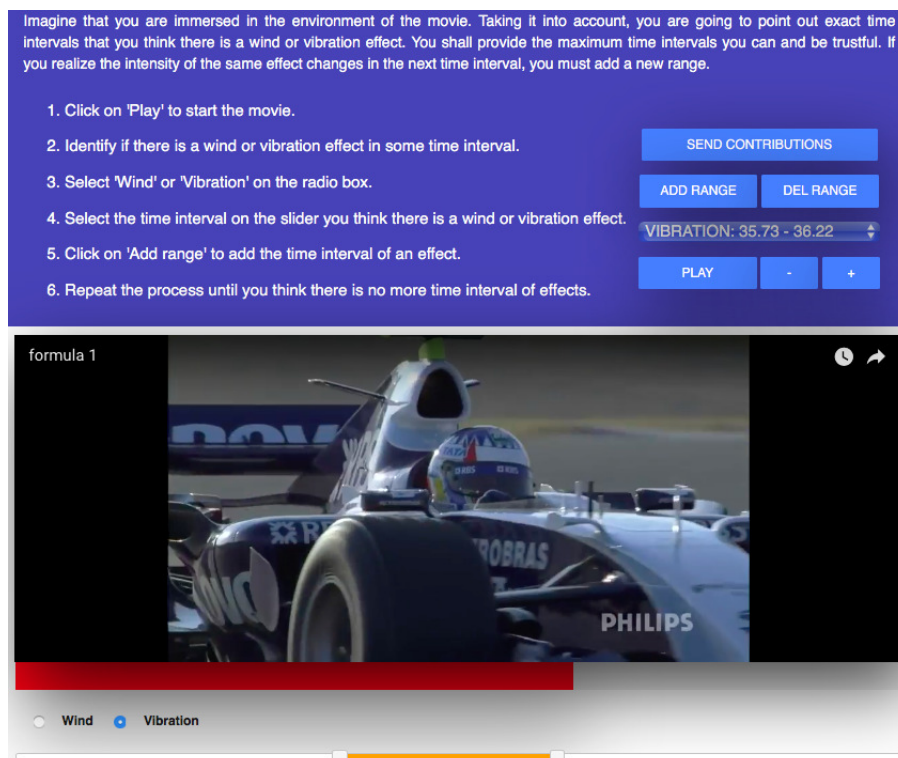


Fig. 5: Tool for identifying sensory effects.

The second microtask was executed in two stages: In the first one was executed 60 jobs and, later, another 90 were executed, totaling 150 jobs. In each job, a worker annotated one or more time spans in which he/she believed that the effects of *Wind* or *Vibration* should be inserted. The output was a list of time ranges of *Wind* and *Vibration* effects. The 150 contributions provided 166 ranges for insertion of *Wind* effects and 146 of *Vibration* effects.

### 5.5.2 Filtering

In an attempt to eliminate malicious and inconsistent contributions, two filtering criteria were used: (i) amount of ranges provided in the contribution, and (ii) existence of overlap between the ranges provided in the contribution.

To meet the first criterion, the average number of ranges in a same contribution for each segment was calculated, and contributions that received less than 50% or more than 200% of that amount were eliminated. The second criterion evaluated the existence of overlap between ranges provided by the same worker for a given effect, in which case it was assumed that the range was updated but the worker forgot to delete the first annotation, so only the most recent ranges of each overlap was maintained .

At the end of the filtering process, 131 of the 146 *Wind* effects identifications and 147 of the 166 *Vibration* ones remained. For the video *Babylon A.D.*, there were 25 identifications of the *Wind* and 46 of the *Vibration* effects; for *Earth* video, 57 identifications of the *Wind* and 46 of the *Vibration* effects and, finally, for *Formula 1* video, 49 identifications of the *Wind* and 55 of the *Vibration* effects.

### 5.5.3 Aggregation

The aggregation of the contributions collected in this second microtask is based on grouping the contributions so that each group contains suggestions from the crowd regarding the same time range for adding a sensory effect.

Firstly, the intervals were divided by video and subdivided by type of effect, *Wind* or *Vibration*. Then each division is ordered in relation to the initial benefit of the range. Finally, a grouping of the intervals is performed so that each group is composed of overlapping ranges. Non-overlapping ranges were discarded. Each of these groups represents a range for insertion of sensory effect, in which there was agreement of the crowd. In this way, the aggregation function was applied to each group to determine the convergent ranges.

The aggregation function determines the maximum degree of overlap between contributions. Then, this maximum degree is used as a boundary to adjust the upper and lower limits of each range in order to delimit the wider region with degree of overlap greater than half of the maximum.

Table 4 shows the numbers of ranges for *Wind* and *Vibration* effects provided by the crowd for each video and its processing. The 150 contributions collected from the workers provided 312 notes of sensory effects, being 146 of *Wind* and 166 of *Vibration*. After filtering, only 278 of 312 were carried forward to the aggregation stage, which in turn delivers 29 *Wind* and 31 *Vibration* effects to be annotated on the selected videos.

Table 4: Contributions and results for the second task.

Video	Wind			Vibration		
	Ranges	Filtered	Converged	Ranges	Filtered	Converged
Babylon A.D.	28	25	8	48	46	10
Earth	65	57	11	60	46	5
Formula 1	53	49	10	58	55	16
<b>Total</b>	<b>146</b>	<b>131</b>	<b>29</b>	<b>166</b>	<b>147</b>	<b>31</b>

## 5.6 Conclusion Stage

The conclusion stage aims to generate the crowdsourcing project outcome. At this point, there already exists the internal representation of the sensory effect metadata for each selected video. To promote interoperability, the final result of the process is represented in conformance with MPEG-V format, so that the

results generated through this work can be refined with the help of tools like SEVino [37] and Real 4D studio [33], and rendered by multimedia players such as PlaySEM [32] and SEMP [37]. The results could also be represented in the format EAF (ELAN Annotation Format), compatible with the ELAN<sup>5</sup> video annotation system which, in addition to allowing the result to be displayed, can also export it to other formats.

## 6 Results and Discussion

In order to analyse the results of our study, we made comparisons of the content produced by our approach, using the three videos selected from the database of Waltl et al. [39] (i.e. *Babylon A.D.*, *Earth* and *Formula 1*), with the annotations for the same three videos, produced by a specialized team responsible for populating this public database. Although this public database contains information about the corresponding sensory effects and their attributes, as intensity, we decided that crowd members should only determine the corresponding, *Wind* or *Vibration*, sensory effect to each video scene, no matter the intensity of the annotated effect.

It is noteworthy that the comparison between the effects identified by the crowd and those identified by the author, rather than measuring the similarity between the results, aims to understand how they complement each other.

### 6.1 Babylon A.D.

The video is a commercial trailer for an action movie that features mainly gunshots and explosions. The workers contributed 38 times to it, indicating 28 *Wind* and 48 *Vibration* effects to the video. After running the task aggregation method, 8 time intervals containing *Wind* and 10 time intervals containing *Vibration* remained. The most noticeable events on this video corresponding to gunshots and explosions were identified by the crowd, including some which hadn't previously been annotated in the reference dataset. Moreover although the two most perceptible explosion events in the reference dataset were not obtained using the aggregation method, these events received contributions from the crowd participants.

Table 5 and Table 6 show the effects of *Wind* and *Vibration* obtained from the crowd compared to the annotations of the reference dataset [39].

While analyzing the content of the *Babylon A.D.* trailer, it was possible to identify why some effects were present in the dataset and not perceived in the same way by the crowd. For instance, at the beginning of the video presents an object (satellite) moving through the space. Although the dataset metadata had a *Wind* effect annotated, the workers did not indicate that probably because they considered that there is no air flow in the space. This occurrence demonstrates how the author tends to use the sensory effects to convey his

---

<sup>5</sup> ELAN available at <https://tla.mpi.nl/tools/tla-tools/elan>

Table 5: Babylon A.D. - Vibration.

Author		Crowd	
start	end	start	end
		8.09	10.99
26.7	30.7	29.04	30.94
		32.85	33.55
		39.28	39.98
		47.02	47.32
49.2	49.6	48.60	57.90
61.6	63.6	61.30	67.00
		69.40	73.00
74.7	74.9		
75.5	75.8		
78.0	78.5		
89.6	99.2	98.04	98.24
99.2	109.2		
		112.00	117.20

Table 6: Babylon A.D. - Wind.

Author		Crowd	
start	end	start	end
10.4	11.8		
12.4	33.5	29.31	33.41
		36.00	39.20
40.0	44.3	42.07	45.97
45.9	49.2	47.00	47.10
53.8	55.0	54.09	57.99
63.1	63.6		
73.8	74.6		
75.5	75.9		
		92.47	93.07
97.2	99.0		
102.6	109.2	105.70	106.90
		113.08	114.88

artistic vision of the scene. On top of that, the crowd tended to associate the effects as they could feel it, like in explosions, whereas the reference dataset associated more effects in scenes with low motion.

By analyzing the annotated video it is possible to verify that the most evident events were identified by the author and the crowd. Gun shots and explosions of lesser intensity were noticed only by the crowd, while the author’s explicit annotations refer to subtle events. In this way, the effects identified by the author and by the crowd are complementary, covering both the workers’ expectations and the subtleties intended by the author.

## 6.2 Earth

The *Earth* trailer had 63 contributions, resulting in a total of 65 *Wind* effects and 60 *Vibration* effects. After running our aggregation method, 11 time intervals containing *Wind* and 5 time intervals containing *Vibration* remained. The workers noticed more *Vibration* effects in this video than the public dataset. Taking into account the analysis of the *Earth* scenes, we concluded that the workers perceived *Vibration* in scenes with stronger movements, such as when an animal jumps or in a herd of animals running. Moreover, when they heard a very loud noise in scene transitions, they pointed *Vibrations*. Tables 7 and 8 show the effects of *Wind* and *Vibration* obtained from the contributions of the crowd, compared to the annotations present in the reference dataset [39] with an intensity greater than or equal to 50%.

Regarding the *Wind* effect, the workers associated it with fast movements, e.g. scenes with cloud movements and a herd of animals running. They also noticed *Wind* in a scene where a quick presentation of slides with animal images was displayed.

It was evident to this video that the crowd complemented the effects indicated by the author, adding *Wind* effects to scenes that featured fast movements and *Vibration* effects for scenes with strong movements.

Table 7: Earth - Vibration.

Author		Crowd	
start	end	start	end
		3.40	6.30
18.3	18.7		
		21.0	27.70
		37.40	39.20
		52.90	53.70
		57.25	57.75

Table 8: Earth - Wind.

Author		Crowd	
start	end	start	end
6.5	10.0	9.97	11.67
12.7	14.0	12.59	12.99
17.7	21.0	18.00	22.90
23.0	29.2	25.00	27.90
33.0	33.9	32.90	33.20
35.6	39.2	34.78	35.78
		38.40	39.10
41.2	44.8		
		47.35	47.45
		55.10	55.90
		59.04	60.94
		63.00	65.00

### 6.3 Formula 1

This AV content is an advertisement for *Formula 1* racing, in which there are several scenes of racing cars as well as scenes of pit stops and pilots walking. The workers made 49 contributions, resulting in a total of 53 *Wind* effects and 58 *Vibration* effects. After running our aggregation method, 10 time intervals containing *Wind* and 16 containing *Vibration* remained. Table 9 and Table 10 show the effects of *Wind* and *Vibration* obtained from the crowd contributions compared to the sensory effect annotations with an intensity greater than or equal to 50% found in our reference dataset.

Similarly, as occurred in the *Babylon A.D.* and *Earth* videos, workers did not notice events associated with the reference dataset to low-intensity tactile stimulus (*Wind* and *Vibration*) in *Formula 1* video. However, they noticed most of the time spans with the intensity stronger than 50%, which may indicate that less pronounced effects are more related to the expression of authorship than something highly expected by most viewers. At the beginning of the clip, drivers were slow and there was an indication of *Wind* in the public dataset whereas the workers did not judge it as they would feel the *Wind* in the environment of the movie. Furthermore, the workers spontaneously indicated *Vibration* when the cars accelerated, which there was not present in the reference dataset. Besides, the crowd covered almost all *Vibration* effects annotated in this dataset even with low intensities.

The analysis of the results for this video shows that, predominantly, the crowd complemented the effects annotated by the author with *Vibration* effects for the events of cars acceleration, and with *Wind* effects for the overtaking events in a racing.

Table 9: Formula 1 - Vibration.

Author		Crowd	
start	end	start	end
7.67	7.97		
		11.60	13.40
		16.00	18.90
		23.60	28.20
31.0	32.3	31.51	31.91
		33.59	36.29
		36.68	37.58
		43.82	44.02
47.0	47.9		
		54.03	56.83
		62.00	66.00
		71.16	71.76
		78.75	78.95
		88.95	91.65
		92.43	96.53
96.7	99.0	98.21	101.21
100.7	101.2		
		102.01	102.71

Table 10: Formula 1 - Wind.

Author		Crowd	
start	end	start	end
9.0	17.0	16.00	18.90
26.6	29.2	28.10	29.20
31.0	32.3		
		34.00	34.20
41.7	43.4		
		45.40	46.80
47.0	47.9		
		52.00	59.50
		63.00	63.10
65.5	66.5		
		69.00	69.20
70.0	72.0		
75.5	79.0		
80.5	91.0	82.30	83.30
96.7	105.0	100.43	100.93
108.0	116.0	108.78	112.38

#### 6.4 Crowdsourcing and mulsemmedia content authoring

At the current stage, the experimental results show that our crowdsourcing approach for sensory effect authoring is a viable alternative when it is combined with an expert fine-tuning of mulsemmedia authoring. While watching the annotated videos, we realized that most of the differences between the MPEG-V SEM from the dataset and from the crowd could be justified. We believe that individuals are subject to natural bias and oversight when authoring mulsemmedia content due to their unique prior experiences and the expected QoE of users is subject to what they are feeling.

A collective common sense to indicate time intervals of sensory effects can be an effective starting point for mulsemmedia content annotation. On the other hand, it is still necessary to incorporate expert advice to fine-tune sensory effects attributes such as intensity, location, and so forth. This approach could be opportune for the industry to turn their multimedia videos into mulsemmedia ones, outsourcing the hard work of pointing out the sensory effects presented in their library, and then, fine-tuning the work with their own experts.

Table 11 summarizes the number of contributions collected in each task, as well as the number of contributions remaining after the application of the

filter criteria and the number of results produced by the aggregation method in each task.

Table 11: Contributions collected and aggregated.

Task	Items	Contributions	Annotations	Filtered	Aggregated
1	3 videos	61	434	306	38 segments
2	38 segments	150	312	278	60 effects
<b>Total</b>	<b>41 items</b>	<b>211</b>	<b>746</b>	<b>584</b>	<b>98 items</b>

In total, 211 contributions were collected during the entire process. Of these contributions, 746 annotations were obtained. Applying the reliability filters during the process, 584 annotations were considered valid, that is, there was a 78.28% effectiveness rate in the contributions.

Including crowdsourcing platform costs, the total spent on the campaign was 12.53 USD. As the amount paid to workers for each job was 0.05 USD, the total cost was 10.55 USD.

Since 29 *Wind* and 31 *Vibration* effects were obtained, the average was 0.43 USD per *Wind* effect and 0.40 USD per *Vibration* effect. These values are summarized in Table 12.

Table 12: Cost of the campaign.

Cost (USD)	
Total	12.53
Contributions	10.55
Platform Fees	1.98
per Wind Effect	0.43
per Vibration Effect	0.42
per Contribution	0.05

## 7 Conclusion

This paper introduces a completely new approach for authoring mulsemmedia content based on crowdsourcing contributions. We contrasted our results to a public mulsemmedia dataset to assess the proximity of the information provided by the crowd. The results pointed that the authoring made by the crowd adds to the public dataset and vice-versa.

A major observation made about the results is that the effects identified by the crowd are largely not the same as those annotated by the author. This was already expected, since the effects obtained by the crowd reflect the workers' point of view regarding the intervals in which they believe that it makes sense to have sensory effects, whereas the author has a greater concern in transmitting his artistic point of view through of effects. The crowd was

able to indicate the semantic associations related to the effects of *Wind* and *Vibration*, however, it was clear that the proposed insertions were for evident effects. The practical use of this is to adopt a hybrid approach in which the author divides the authoring with the multitude, delegating to the workers the work of identifying the insertions of greater intensity that relate to the user experience, allowing the author to concentrate on authoring more refined representation of his artistic vision. In other words, the authoring of the crowd did not replace the author's, but rather the help becomes more appropriate to improve the quality of the user experience.

For instance, in *Formula 1* the workers almost always realized *Wind* effects and *Vibration* when cars were accelerating the cars. These effects were not authored in the reference dataset. This draws attention to the possibility of using the proposed approach to complement annotations made by experts.

Equally important, the idea behind our approach relies on the combination of the intuitive judgment of several individuals, the common sense, and the refinement of an expert in order to take the best of each perspective to provide an alternative method for authoring mulsemmedia content. As defined by van Holthoorn and Olson [19] "*common sense consists of knowledge, judgment, and taste which is more or less universal and which is held more or less without reflection or argument.*" The presented approach takes advantage of common sense emerged from the crowd in terms of expected sensory effects associated with each video scene. Also, it could be timely for mass production of coherent mulsemmedia content without taking endless hours of an expert to start from the scratch.

It is worth noticing that mulsemmedia annotation does not automatically lead to mulsemmedia authoring. For instance, in the case of *Wind* effects, where there is a lingering effect, the fact that the crowd hasn't identified the *Babylon A.D.*'s segment [73.8, 74.6] seconds does not necessarily mean that the fan has to be switched on/off at these points but it is a cue. Indeed, because of lingering effects, network and device delays, a delta for propagation delay should be considered.

Another important observation is that even with a limited number of contributions the crowd associated sensory effects for the most evident situations, such as when explosions, gunshots or car accelerations occur in scenes. Hence, this approach could be applied to build larger datasets storing video annotated with sensory effects. Moreover, these datasets could be used for training systems based on machine learning to detect the previous situations in AV contents automatically. In addition, our approach can also be used for QoE evaluation purposes such as the work of Yue, Wang and Cheng [41].

Future work includes finding how to use the wisdom of the crowd to collect fine-tuned attributes to the maximum extent as well as the automatic generation of MPEG-V metadata for mulsemmedia content without needing traditional annotation tools to go with the process. Furthermore, similarly to this work, the annotation of other types of effects such as olfactory and thermal, using the cascade crowdsourcing process, might be included in next experiments.



**Acknowledgements** Estêvão Bissoli Saleme thankfully acknowledges support from the Federal Institute of Espírito Santo. Prof. Gheorghita Ghinea gratefully acknowledges funding from the European Union's Horizon 2020 Research and Innovation programme under Grant Agreement no. 688503 for the NEWTON project (<http://www.newtonproject.eu>). The authors would also like to thank FAPES, CAPES, and CNPq for financial support for this research.

## References

1. Ademoye, O.A., Murray, N., Muntean, G.M., Ghinea, G.: Audio masking effect on inter-component skews in olfaction-enhanced multimedia presentations. *ACM Trans. Multimedia Comput. Commun. Appl.* **12**(4), 51:1–51:14 (2016). DOI 10.1145/2957753
2. Amorim, M.N., Neto, F.R.A., Santos, C.A.S.: Achieving complex media annotation through collective wisdom and effort from the crowd. In: 2018 25th International Conference on Systems, Signals and Image Processing (IWSSIP), pp. 1–5. IEEE (2018). DOI 10.1109/IWSSIP.2018.8439402
3. Ballan, L., Bertini, M., Del Bimbo, A., Seidenari, L., Serra, G.: Event detection and recognition for semantic annotation of video. *Multimedia Tools and Applications* **51**(1), 279–302 (2011). DOI 10.1007/s11042-010-0643-7. URL <https://doi.org/10.1007/s11042-010-0643-7>
4. Bartocci, S., Betti, S., Marcone, G., Tabacchiera, M., Zanucoli, F., Chiari, A.: A novel multimedia-multisensorial 4d platform. In: AEIT International Annual Conference (AEIT), 2015, pp. 1–6. IEEE (2015). DOI 10.1109/AEIT.2015.7415215
5. Chen, J., Yao, T., Chao, H.: See and chat: automatically generating viewer-level comments on images. *MTAP: Multimedia Tools and Applications* pp. 1–14 (2018). DOI <https://doi.org/10.1007/s11042-018-5746-6>
6. Cho, H.: Event-based control of 4d effects using mpeg rose. Master's thesis, School of Mechanical, Aerospace and Systems Engineering, Division of Mechanical Engineering. Korea Advanced Institute of Science and Technology. Master's Thesis (2010)
7. Choi, B., Lee, E.S., Yoon, K.: Streaming media with sensory effect. In: Information Science and Applications (ICISA), 2011 International Conference on, pp. 1–6. IEEE (2011). DOI 10.1109/ICISA.2011.5772390
8. Chowdhury, S.N., Tandon, N., Weikum, G.: Know2look: Commonsense knowledge for visual search. In: Proceedings of the 5th Workshop on Automated Knowledge Base Construction, pp. 57–62 (2016)
9. Covaci, A., Zou, L., Tal, I., Muntean, G.M., Ghinea, G.: Is multimedia multisensorial?-a review of mulsemmedia systems. *ACM Computing Surveys (CSUR)* **51**(5), 91 (2018)
10. Cross, A., Bayyapunedi, M., Ravindran, D., Cutrell, E., Thies, W.: Vidwiki: enabling the crowd to improve the legibility of online educational videos. In: Proceedings of the 17th ACM conference on Computer supported cooperative work & social computing, pp. 1167–1175. ACM (2014)
11. Di Salvo, R., Spampinato, C., Giordano, D.: Generating reliable video annotations by exploiting the crowd. In: IEEE Winter Conf. on Applications of Computer Vision (WACV), pp. 1–8. IEEE (2016). DOI 10.1109/WACV.2016.7477718
12. Dumitrache, A., Aroyo, L., Welty, C., Sips, R.J., Levas, A.: "dr. detective": Combining gamification techniques and crowdsourcing to create a gold standard in medical text. pp. 16–31 (2013)
13. Egan, D., Brennan, S., Barrett, J., Qiao, Y., Timmerer, C., Murray, N.: An evaluation of heart rate and electrodermal activity as an objective qoe evaluation method for immersive virtual reality environments. In: Eighth International Conference on Quality of Multimedia Experience (QoMEX'16) (2016). DOI 10.1109/QoMEX.2016.7498964
14. Foncubierta Rodríguez, A., Müller, H.: Ground truth generation in medical imaging: A crowdsourcing-based iterative approach. In: Proceedings of the ACM Multimedia 2012 Workshop on Crowdsourcing for Multimedia, CrowdMM '12, pp. 9–14. ACM, New York, NY, USA (2012). DOI 10.1145/2390803.2390808
15. Galton, F.: Vox populi (the wisdom of crowds). *Nature* **75**(7), 450–451 (1907)

16. Ghinea, G., Timmerer, C., Lin, W., Gulliver, S.R.: Mulsemmedia: State of the art, perspectives, and challenges. *ACM Trans. Multimedia Comput. Commun. Appl.* **11**(1s), 17:1–17:23 (2014). DOI 10.1145/2617994
17. Gottlieb, L., Choi, J., Kelm, P., Sikora, T., Friedland, G.: Pushing the limits of mechanical turk: qualifying the crowd for video geo-location. In: *Proceedings of the ACM multimedia 2012 workshop on Crowdsourcing for multimedia*, pp. 23–28. ACM (2012)
18. Hardman, L., Obrenović, Ž., Nack, F., Kerhervé, B., Piersol, K.: Canonical processes of semantically annotated media production. *Multimedia Systems* **14**(6), 327–340 (2008). DOI 10.1007/s00530-008-0134-0. URL <https://doi.org/10.1007/s00530-008-0134-0>
19. van Holthoon, F., Olson, D.: *Common Sense: The Foundations for Social Science*. Common Sense. University Press of America (1987)
20. Kim, S., Han, J.: Text of white paper on mpeg-v. Tech. Rep. ISO/IEC JTC 1/SC 29/WG 11 W14187, San Jose, USA (2014)
21. Kim, S.K.: Authoring multisensorial content. *Signal Processing: Image Communication* **28**(2), 162–167 (2013). DOI 10.1016/j.image.2012.10.011
22. Kim, S.K., Yang, S.J., Ahn, C.H., Joo, Y.S.: Sensorial information extraction and mapping to generate temperature sensory effects. *ETRI Journal* **36**(2), 224–231 (2014). DOI 10.4218/etrij.14.2113.0065
23. Lasecki, W., Miller, C., Sadilek, A., Abumoussa, A., Borrello, D., Kushalnagar, R., Bigham, J.: Real-time captioning by groups of non-experts. In: *Proceedings of the 25th annual ACM symposium on User interface software and technology - UIST '12, UIST '12*, pp. 23–33. ACM Press, New York, New York, USA (2012). DOI 10.1145/2380116.2380122
24. Masiar, A., Simko, J.: Short video metadata acquisition game. In: *10th International Workshop on Semantic and Social Media Adaptation and Personalization (SMAP)*, pp. 61–65. IEEE (2015). DOI 10.1109/SMAP.2015.7370092
25. McNaney, R., Othman, M., Richardson, D., Dunphy, P., Amaral, T., Miller, N., Stringer, H., Olivier, P., Vines, J.: Speeching: Mobile crowdsourced speech assessment to support self-monitoring and management for people with parkinson’s. In: *Proc. of the 2016 CHI Conf. on Human Factors in Computing Sys. - CHI '16, CHI '16*, pp. 4464–4476. ACM Press, New York, New York, USA (2016). DOI 10.1145/2858036.2858321
26. Murray, N., Lee, B., Qiao, Y., Muntean, G.M.: The influence of human factors on olfaction based mulsemmedia quality of experience (2016). DOI 10.1109/QoMEX.2016.7498975
27. Neto, F.R.A., Santos, C.A.S.: Understanding crowdsourcing projects: A systematic review of tendencies, workflow, and quality management. *Information Processing & Management* **54**(4), 490–506 (2018). DOI <https://doi.org/10.1016/j.ipm.2018.03.006>
28. Oh, H.W., Huh, J.D.: Auto generation system of mpeg-v motion sensory effects based on media scene. In: *Consumer Electronics (ICCE), 2017 IEEE International Conference on*, pp. 160–163. IEEE (2017). DOI 10.1109/ICCE.2017.7889269
29. Rainer, B., Waltl, M., Cheng, E., Shujau, M., Timmerer, C., Davis, S., Burnett, I., Ritz, C., Hellwagner, H.: Investigating the impact of sensory effects on the quality of experience and emotional response in web videos. In: *Fourth International Workshop on Quality of Multimedia Experience (QoMEX)*, pp. 278–283. IEEE (2012). DOI 10.1109/QoMEX.2012.6263842
30. Sadallah, M., Aubert, O., Prié, Y.: Chm: an annotation- and component-based hyper-video model for the web. *Multimedia Tools and Applications* **70**(2), 869–903 (2014). DOI 10.1007/s11042-012-1177-y. URL <https://doi.org/10.1007/s11042-012-1177-y>
31. Saleme, E.B., Celestrini, J.R., Santos, C.A.S.: Time evaluation for the integration of a gestural interactive application with a distributed mulsemmedia platform. In: *Proceedings of the 8th ACM on Multimedia Systems Conference, MMSys'17*, pp. 308–314. ACM, ACM, New York, NY, USA (2017). DOI 10.1145/3083187.3084013
32. Saleme, E.B., Santos, C.A.S., Ghinea, G.: Coping with the challenges of delivering multiple sensorial media. *IEEE MultiMedia* pp. 1–1 (2018). DOI 10.1109/MMUL.2018.2873565
33. Shin, S.H., Ha, K.S., Yun, H.O., Nam, Y.S.: Realistic media authoring tool based on mpeg-v international standard. In: *Ubiquitous and Future Networks (ICUFN), 2016 Eighth International Conference on*, pp. 730–732. IEEE (2016). DOI 10.1109/ICUFN.2016.7537133

34. Taborsky, E., Allen, K., Blanton, A., Jain, A.K., Klare, B.F.: Annotating unconstrained face imagery: A scalable approach. In: *Inter. Conf. on Biometrics (ICB)*, pp. 264–271. IEEE (2015). DOI [10.1109/ICB.2015.7139094](https://doi.org/10.1109/ICB.2015.7139094)
35. Teki, S., Kumar, S., Griffiths, T.D.: Large-scale analysis of auditory segregation behavior crowdsourced via a smartphone app. *PLoS ONE* **11**(4) (2016). DOI [10.1371/journal.pone.015](https://doi.org/10.1371/journal.pone.015)
36. Timmerer, C., Waltl, M., Rainer, B., Hellwagner, H.: Assessing the quality of sensory experience for multimedia presentations. *Signal Processing: Image Communication* **27**(8), 909–916 (2012). DOI <https://doi.org/10.1016/j.image.2012.01.016>
37. Waltl, M., Rainer, B., Timmerer, C., Hellwagner, H.: An end-to-end tool chain for sensory experience based on mpeg-v. *Signal Processing: Image Communication* **28**(2), 136–150 (2013). DOI [10.1016/j.image.2012.10.009](https://doi.org/10.1016/j.image.2012.10.009)
38. Waltl, M., Timmerer, C., Hellwagner, H.: Improving the quality of multimedia experience through sensory effects. In: *Second International Workshop on Quality of Multimedia Experience (QoMEX)*, pp. 124–129. IEEE (2010)
39. Waltl, M., Timmerer, C., Rainer, B., Hellwagner, H.: Sensory effect dataset and test setups. In: *Fourth International Workshop on Quality of Multimedia Experience (QoMEX)*, pp. 115–120. IEEE (2012). DOI [10.1109/QoMEX.2012.6263841](https://doi.org/10.1109/QoMEX.2012.6263841)
40. Yuan, Z., Bi, T., Muntean, G.M., Ghinea, G.: Perceived synchronization of multimedia services. *IEEE Transactions on Multimedia* **17**(7), 957–966 (2015). DOI [10.1109/TMM.2015.2431915](https://doi.org/10.1109/TMM.2015.2431915)
41. Yue, T., Wang, H., Cheng, S.: Learning from users: a data-driven method of qoe evaluation for internet video. *MTAP: Multimedia Tools and Applications* pp. 1–32 (2018). DOI <https://doi.org/10.1007/s11042-018-5918-4>
42. Zhai, H., Lingren, T., Deleger, L., Li, Q., Kaiser, M., Stoutenborough, L., Solti, I.: Web 2.0-based crowdsourcing for high-quality gold standard development in clinical natural language processing. *Journal of Medical Internet Research* **15**(4), 1–17 (2013). DOI [10.2196/jmir.2426](https://doi.org/10.2196/jmir.2426)