

Mulsemmedia: State-of-the- Art, Perspectives and Challenges

GHEORGHITA GHINEA Brunel University
CHRISTIAN TIMMERER Alpen-Adria-Universität
WEISI LIN Nanyang Technological University
STEPHEN R. GULLIVER University of Reading

Mulsemmedia – multiple sensorial media – captures a wide variety of research efforts and applications. This paper presents a historic perspective on mulsemmedia work and reviews current developments in the area. These take place across the traditional multimedia spectrum – from virtual reality applications to computer games - as well as efforts in the arts, gastronomy and therapy, to mention a few. We also describe standardization efforts, via the MPEG-V standard, and identify future developments and exciting challenges the community needs to overcome.

Categories and Subject Descriptors: **H.5.2 [Information Interfaces and Presentation]:** User Interfaces—*Evaluation/methodology*; **H.1.2 [Models and Principles]:** User/Machine Systems—*Human Information Processing*;

General Terms: Mulsemmedia, multi-sensory

Additional Key Words and Phrases: Contour perception, flow visualization, perceptual theory, visual cortex, visualization

ACM Reference Format:

Ghinea G, Timmerer, C., Lin, W. and Gulliver, S.R.. Mulsemmedia: State-of-the- Art, Perspectives and Challenges. ACM Trans. Multimedia Computing Communications and Applications. X, Y, Article Z (XXX 201X), XX pages.

DOI=10.1145/0000000.0000000 <http://doi.acm.org/10.1145/0000000.0000000>

1. INTRODUCTION

In 2004, on the 10th anniversary of the creation of the ACM Multimedia Special Interest Group, Larry Rowe and Ramesh Jain published a seminal paper on “*Future Directions in Multimedia Research*” in ACM TOMCCAP. The paper presented the result of discussions at a one day workshop with over 30 leading researchers in the field. There was agreement that multimedia is a multidisciplinary field, applying to a variety of fields (e.g., entertainment, education, medicine, creative arts, etc.). Three unifying themes were identified to unite the multimedia research field. Firstly, a multimedia system or application is comprised of at least two media objects that are correlated. Secondly, there is the issue of integration and adaptation where multiple media objects should be used jointly and separately to improve application performance, and distributed multimedia applications should provide transparent delivery of dynamic content in such a way that content adapts naturally to the users’ environment. Thirdly, multimedia applications are multimodal and interactive (Rowe and Jain, 2005).

10 years on, what has changed? A lot, and maybe not so. Arguably, the three unifying themes are very much valid today, in a world dominated by social media and a proliferation of sensor rich (predominantly mobile) devices, where individuals are producers, broadcasters, and consumers of rich media content. Reassuringly, the accepted definition of multimedia remains that of a combination of two or more media, one of which is preferably continuous, the other usually discrete. It is without doubt that most of multimedia content available today is a combination of video and audio (both continuous media) with textual (discrete media) information sometimes contained therein. However, such applications engage primarily two of our human senses: that of sight and hearing, i.e. they are

Author’s address: G. Ghinea, Department of Computer Science, Kingston Lane, Uxbridge, UB8 3PH, UK; email: george.ghinea@brunel.ac.uk; Christian Timmerer, Universitätsstrasse 65-67 A-9020 Klagenfurt Austria; email: christian.timmerer@itec.uni-klu.ac.at; Weisi Lin, School of Computer Engineering, Nanyang Technological University 50 Nanyang Avenue Singapore 639798; email: wslin@ntu.edu.sg; Stephen Gulliver, Henley Business School, Whiteknights, Reading, RG6 6UR, UK; email: s.r.gulliver@henley.reading.ac.uk

Permission to make digital or hardcopies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies show this notice on the first page or initial screen of a display along with the full citation. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credits permitted. To copy otherwise, to republish, to post on servers, to redistribute to lists, or to use any component of this work in other works requires prior specific permission and/or a fee. Permissions may be requested from Publications Dept., ACM, Inc., 2 Penn Plaza, Suite 701, New York, NY 10121-0701 USA, fax +1 (212) 869-0481, or permissions@acm.org.

@2010 ACM 1544-3558/2010/05-ART1 \$10.00

DOI10.1145/0000000.0000000 <http://doi.acm.org/10.1145/0000000.0000000>

bi-sensorial. This situation is at odds with the fact that 60% of human communication is non-verbal and that most of us perceive the world through a combination of five senses (i.e., sight, hearing, touch, taste, and smell). As such, current multimedia experiences fail to convey the sensation, for instance, of heat and humidity, let alone the wafts of aromas that one experiences when waking through a spice market in India. As humans, we engage and learn by interacting with all of our senses – can we not do this in a digital fashion as well?

We therefore propose Mulsemmedia – multiple sensorial media - as a new multimedia challenge for the forthcoming 10 years. Whereas multimedia applications are usually bi-(sometimes tri-)media and almost exclusively bi-sensorial in nature, mulsemmedia applications are those that engage three (or more) of our senses.

While current technological developments have made *digital* mulsemmedia experiences somewhat of a novelty, in the non-digital world they are anything but. The earliest we know about happened in 1906 when artificially generated smells were combined with audiovisual content. An audience was sprayed with the scent of roses while watching a screening of the Rose Bowl football game. In 1943, Hans Laube who had earlier perfected a technique to extract odors from an enclosed environment, was able to reverse this process so that selected scents were emitted at specific times and for specified durations, resulting in a 35 minute ‘smell-o-drama’ movie called *Mein Traum* in which 35 different odors were released to accompany the drama presentation. Building on this, audiences in 1959 viewing a documentary about Red China called *Behind the Great Wall* were treated with an AromaRama presentation, in which the theatre’s air-conditioning system was used to release over 30 different smells. Shortly afterwards, in 1960, Michael Todd Jr produced a competing system called *Smell-O-Vision*, in which aromas were released during the screening of the movie Scent of Mystery. It would be an exaggeration to say that these experiences were an unqualified success: challenges of generating realistic scents, the tendency of odors to drift and diffuse, as well as insufficiently understood characteristics of odor intensity all meant that, novelty factor aside, user take-up was low. The reaction of the audience to the AromaRama experience is probably best described from the following extract from the review published back then by Time magazine:

“To begin with, most of the production’s 31 odors will probably seem phoney, even to the average uneducated nose. A beautiful old pine grove in Peking, for instance, smells rather like a subway rest room on disinfectant day. Besides, the odors are strong enough to give a bloodhound a headache. What is more, the smells are not always removed as rapidly as the scene requires: at one point, the audience distinctly smells grass in the middle of the Gobi desert.”

Such drawbacks did not prevent pioneering mulsemmedia efforts, however. In 1962, Morton Heilig created what is now popularly dubbed as the first virtual reality (VR) experience for users, even though digital computing and virtual reality systems did not exist then. With *Sensorama*, he created an arcade-style device, which took users on an immersive 3-D virtual reality bike ride experience through the streets of Brooklyn, New York. This came complete with motions and vibrations, sounds, fans and smells, the most complex mulsemmedia experience devised so far, engaging four out of our five major senses. Indeed, given that the sense of taste is intimately connected to that of smell and that one of the aromas emitted was that of freshly baked bread from a bakery, it is not inconceivable that for some users all five major senses were engaged in their mulsemmedia journey (Heilig, 1962).

Over half a century has passed since then – so where are we now on the mulsemmedia landscape? To answer the question, this paper reviews developments that recent technological advances have made possible to see how mulsemmedia applications fit within the multimedia arena, and to identify challenges that the community has to overcome. Accordingly, the structure of the rest of this paper is as follows: given the importance of the human sense to mulsemmedia, the next section gives an overview; Section 3 then details related work. Mulsemmedia needs standards to thrive, and, to this end, Section 4 describes MPEG-V a standard capable of supporting mulsemmedia applications. The user is an important element of mulsemmedia, and QoE efforts in this respect are detailed in Section 5. Finally, research challenges and open issues are described in Section 6.

2. HUMAN SENSORIAL OVERVIEW

In this section we consider in more detail the multiple process steps required to achieve multiple sensory perception. We introduce key physiological systems, and describe how each captures and transforms information from the world, so that the brain can process it. We conclude the section by considering the issue of cognitive binding, and highlight the attentive struggle between top-down and bottom-up processes.

2.1 Multiple Sensory Perception

Sensory perception relates to a human's conscious sensory experience of the world, i.e. what a person can see, hear, smell, touch, and taste, etc. When we consider mulsemmedia perception, therefore, it is critical to appreciate that multiple sensory media perception is not something that just 'happens'. For a person to be able to understand and assimilate meaning from multiple sensory media, they must capture, interpret and combine information from numerous sensory organs – bottom up sensing (Goldstein, 2013). Moreover, information from multiple senses must be cognitively joined and aligned, and then compared to higher-order cognitive schema, which define task semantics, pragmatics and social norms – top down thinking (Marois, 2005; Mayer, 2003).

Although perception sometime feels as though it just happens, it is in reality the result of a complex set of processes. Biological sensors capture physical signals from the environment and transduce them, with the exception of specific chemoreceptors, into structured electrical signals. These signals are restructured in the nervous systems, and transmitted to the brain. Within the brain, spatial/temporal signals are then sub-consciously structured as patterns, which are attentively processed as higher-level artefacts / objects. Once structured, appreciation of meaning facilitates the validation of propositions. Identifying whether something is true or false facilitates humans align bottom-up sensory input with top-down knowledge and memory; and enables us to create, and iteratively validate, complex schema models of the real world. The existence of these complex schema models allows humans to predict, and understand, the world in context of higher pragmatic and social structures.

2.2 Capturing the Physical

There is no universal agreement as to the number of senses perceived by the human mind. In reality, however, the human body manages sensory inputs from a wide range of internal and external sensory inputs; such as pain (nociception), space (proprioception), movement (kinaesthesia), time, and temperature (thermoception). As well as external senses, our bodies sense and processes internal regulation (called interoceptive senses), which leads to feelings of hunger, sickness, thirst, stress or discomfort (Craig, 2003). All of these senses are internally linked within our model of the world. Despite our processing this dynamic range of internal and external senses, mulsemmedia systems focus on the five traditional sense, as defined by Aristotle, i.e. visual (sight), auditory (sound), tactile/haptic (touch), Olfactory (smell) and Gustatory (taste). In this section we introduce the reader to each physiological system in turn.

2.2.1 Visual (Sight). In mulsemmedia, sight allows assimilation of textual and visual information. Light reflected from a physical object in the visual field enters the eye through the pupil and passes through the lens; which projects an inverted image onto the retina at the back of the eye. The retina consists of approximately 127 million light-sensitive cells (120 million called rods; 7 million called cones, which can be subdivided into L-cones, M-cones and S-cones). Although cones are less light sensitive than rods, they are responsible for capturing color within the human visual system. When light enters the eye, it passes through seven sensory cell-layers before reaching the rods and cones at the back of the eye. If cones were distributed evenly across the retina, their average distance apart would be relatively large, leading to poor spatial acuity. Accordingly cones are concentrated in the center of the retina (in a circular area called macula lutea). Within this area, there is a depression called the fovea, which consists almost entirely of cones, and it is through this area of high acuity, extending over just 2° of the visual field, that humans make their detailed observations of the world. The cells that process and transmit information to the brain are called the bipolar, horizontal and

ganglion cells. Photoreceptors at the back of the eye (cones and rods) are activated when light is shined at them, which consecutively activates bipolar cells. Visual pre-attentive segregation, and object combining, occurs primarily in the occipital lobe (at the back of the brain), however visual information is contextualized, i.e. 'Where/How' and 'What', in the Parietal and Temporal lobes respectively (Schiller, 1986).

2.2.2 *Auditory (Sound)*. In mulsemmedia, the human auditory system is used heavily in the transfer of sound, speech, music and special effects. If an object vibrates, it produces a sequence of wave compressions in the air surrounding it. These fluctuations in air pressure spread away from the source of vibration at 320m/s, reducing in magnitude as the energy is dispersed. When two or more waveforms interact, they create a combined waveform that is the sum of its component parts. Sound is the sensation produced by the ear when a vibration occurs within a given frequency range (approximately 20 Hz to 20 KHz), which is audible to humans. The volume of sound, at the source of vibration, is dependent upon the magnitude of sound energy waveform. The frequency is dependent upon the frequency of compressions being produced by the source of vibration.

The ear is divided into three parts - the outer (external), the middle and the inner (internal) ear. The outer ear collects sound waves and focuses them along the ear canal to the eardrum. The eardrum vibrates, causing bones (Malleus and Incus) to rock back and forth, which passes movement to the cochlea where fluid in the inner ear is disturbed. The disturbance of fluid causes thousands of small hair cells to vibrate. The cochlea converts sound waves into electrical impulses, which are passed on to the brain via the auditory nerve. The three main auditory areas in the brain (i.e. the core area, the belt area, and the parabelt) are found in the temporal lobe. Recognition of sound and localization of sound are, however, processed separately (Yost, 1985).

2.2.3 *Tactile/ haptic (Touch)*. In mulsemmedia, tactile feedback allows us to identify several distinct types of sensations; as human skin contains a number of different sensory receptor cells that respond preferentially to various mechanical, thermal or chemical stimuli. The majority of multimedia studies involves the tactile or touch sense, which detects pressure and touch (i.e., brushing, vibration, flutter and indentation), however, human skin is also sensitive to temperature and pain. Information from the skin receptors is carried along "touch-neuron pathway" to the somatosensory cortex, which maps the senses in the body and transmits messages about sensory information to other parts of the brain (e.g. for use in performing actions, for making decisions, enjoying sensation or reflecting on them).

2.2.4 *Olfactory (smell)*

In mulsemmedia, olfactory feedback allows researchers to monitor subconscious reaction to smell; which is often linked to task / emotional contextualization. There are 50 million primary sensory receptor cells in a small (2.5 cm²) area of the nasal passage called the olfactory region. The olfactory region is formed of cilia projecting down out of the olfactory epithelium into a layer of mucous, which helps to transfer soluble odorant molecules to the receptor neurons. The neuronal cells form axons, which penetrate the cribriform plate of bone, thus reaching the olfactory bulb of the brain. Smell messages are sent directly to the higher levels of the central nervous system, via the olfactory tract, where olfactory information is decoded and a reaction is determined. Compared to many mammals, smell ability in humans is limited. Smell is, however, important to human perception of episodic knowledge, with smells often triggering specific contextual memories. The olfactory sense is used in humans as a means of identifying resources, as a warning of danger (e.g. rotten food, chemical dangers, and fires), identify mates, predators, aiding navigation, and providing sensual pleasure. Since olfactory neurons are connected directly to the brain, and can therefore unconsciously influence cognition and emotion, smell is known to trigger discomfort, sympathy or even unconscious refusal (Ayabe-kanamura et al., 1998).

2.2.5 *Gastronomy (taste)*.

The tongue is covered in papillae, which are either i) filiform, found across the entire surface of the tongue, ii) fungiform, which are found on the tip and sides of the tongue, iii) foliatae, which are structured at the sides at the back of the tongue, and iv) circumvallate, found at the central back of the tongue. All papillae, with the exception of filiform, contain taste buds. Each of the 10,000 taste

buds contains between 50-100 taste cells. Traditionally it was believed that taste was grouped in areas relating to sour, sweet, salty, bitter (with Umami not considered), however it is now understood that all tastes (including Umami) are registered by all taste buds. Electrical signals are generated in taste buds, and pass along one of a number of nerves, relating to separate areas of the tongue, and link to both the Thalamus (perched on top of the brainstem) and the frontal lobe.

Smell and taste are commonly considered together, as they are functionally linked. Unlike other senses, which interpret light/sound waveforms or interaction patterns, and transform these into electrical signals understood in the brain, smell and taste are often termed ‘gatekeeping’ senses, i.e. sensations created as a result of interaction with molecules being assimilated into the body (Goldstein, 2013). Gatekeeper (chemoreceptor) senses are understandably linked to biological and emotional processes, i.e. to ensure automatic rejection in the case of bad food. Despite input of data via separate sensory systems (i.e., smell and taste), it is almost impossible to taste something while pinching your nose, making the experiences of smell and taste hard to separate.

2.3 Binding and Focus

Although entities and events in the world are perceived via disparate sensory modalities, as described in section 2.2, our experience of the world is largely coherent (both spatially and temporally). The issue of how the brain integrates and aligns sensory fragments is called ‘Binding’ (Damasio, 1989, p. 29), and consists of: segregation and combining processes. Segregation processes (BP1) define high-level object variables within each sensory input (e.g., shape and color from the same input from millions of light-sensitive cells), and combining processes (BP2) relates to the process of joining and synchronizing object variables across different senses. Sensory processing is consistent for all humans, and researchers understand a significant amount concerning the processing and representation of sensory data (Smythies, 1994, p. 54), however there is less understanding of how brain mechanisms construct phenomenal objects (i.e. high-level mental object, either physical or conceptual, which act as the focus of attention).

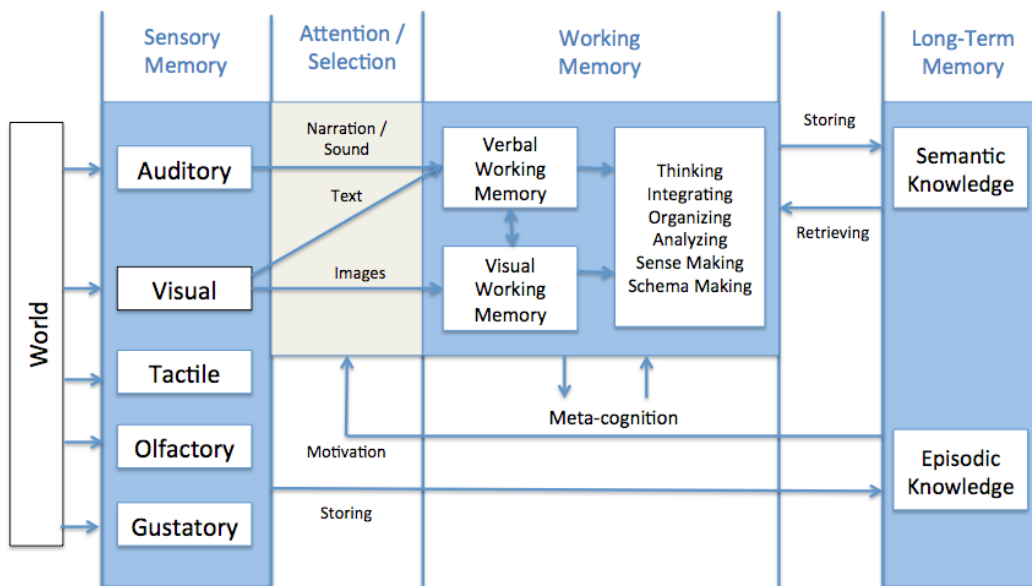


Fig. 1. Schema of Thinking Processes (Based on Mayer, 2003; Marois, 2005; Adapted from Fadel & Lemke, 2008).

Dual process theory, which provides some interesting insights, separates cognition into two systems, i.e. intuition / experience (termed system 1), and reasoning / memory (termed system 2). System 1 combines sensory and emotional stimuli to subconsciously define spatial/temporal associations between object variables. System 2 allows humans to uniquely process conscious judgments and attitudes in context of semantic and episodic knowledge. System 2 is slow, however, and limited in part due its reliance on limited-capacity serial-based memory functions (see Fig 1.).

Such limitations mean that conscious perception occurs in linear installments, with task efficiency significantly reduced if multitask switching is required (Rubinstein et al, 2001). Due to its limitations, human reasoning is significantly dependent upon existing knowledge (schema) to support simplification of the task, or contextualize episodic information; and has been shown to influence user attentive selection (Yarbus, 1967).

2.4 Summary

It is clear that mulsemmedia sensory media perception is not something that just ‘happens.’ Perception is a complex combination of steps that combine bottom-up (sensory processing) and top-down (cognitive reasoning) processes, which result in the appreciation of the media information and the interpretation of its meaning in context of existing semantic and episodic knowledge. Sensory processing is well understood. The process of understanding how knowledge impacts mulsemmedia media interpretation, perception and acceptance, however, is an exciting area of research.

3. RELATED WORK

Mulsemmedia research, while not mainstream and sheltering perhaps under more traditional research areas, has nonetheless progressed over the past 20 years. In this section, we present key work in the area. We start off by highlighting work done on mono-sensorial evaluation. Of course, most work performed so far in this respect targets audition and vision. Since the emphasis of this paper is on mulsemmedia, we will not discuss in detail perception-based models for speech, audio, image, graphics and video; interested readers can refer to the recent surveys in such modeling and applications (e.g., Dermot, et al. 2009, Lin and Kuo 2011, Möller, et al. 2011, Reinhard, et al. 2013, Richard, et al. 2013, Wu, et al. 2013, You, et al. 2010). Therefore, the main thrust of the first sub-section below is to introduce existing research involving other senses, namely, olfaction, taction and gustation – but not in combination with one another. The next subsection then proceeds to review research exploring the combination of two or more senses in a digital environment: mulsemmedia, while Appendix A gives more details for the related basic technical approaches and computational models so far in the literature, although the development of mulsemmedia algorithm and systems is still in its infancy.

3.1 Mono sensorial evaluation

There has been interesting and substantial research into the olfactory system that enables humans to recognize and categorize different odors and determine many behavioral and social reactions. Ho and Spence (2005) investigated the differential effects of olfactory stimulation under conditions of varying task difficulty. Participants detect visually presented target digits from a stream of visually presented distractor letters in a rapid serial visual presentation task; at the same time, participants were required to discriminate stimuli presented on the front or back of their torso. The results showed a significant performance improvement in the presence of peppermint odor (as compared to air) in a difficult task but not in an easy one. This demonstrated that olfactory stimulation can facilitate tactile performance.

In the digital world, a pioneer in the area of olfaction is Kaye (2001), who, in his work on symbolic olfactory devices, experimented with a few prototypical designs of olfactory data display devices to illustrate the concept of computer-controlled smell output. For human beings, odor stimuli are highly associated with many processes such as emotions, attraction, mood, etc. Monitoring and analyzing electroencephalogram (EEG) of human brain activity during perception of odors have shown (Yazdani, et al., 2012) that classification of EEG signals during perception of odors can reveal the pleasantness of the odor with relatively high accuracy. Ghinea, et al. (2010, 2011) focused on olfaction-enhanced applications. The challenges of enhancing mulsemmedia with olfaction were also discussed.

Taction is another important sense for mulsemmedia investigation. For foundational knowledge in this area and guidelines of design, readers can refer to the paper by Seungmoon and Kuchenbecker (2013). Haptic rendering (or haptic display) conveys information about virtual objects to users through the sense of touch. For haptic rendering, force-feedback display of contact interactions can be realized for both rigid and deformable virtual models. A general framework for force-feedback display of virtual environments is presented by Otaduy, et al. (2013), and the issues, modelling, and assessment

related to haptic aesthetics is discussed by Carbon and Jakesch (2013). It has been shown that perceiving material properties (including roughness, friction, and thermal properties) of objects through touch is generally superior to the perception of shape (Klatzky, et al. 2013).

With respect to gustation, this sense is intricately linked with olfaction. However, the only work targeting gustation *per se* of which we are aware is that of Adrian Cheok (<http://adriancheok.info>). He and his team developed a taste transmitter machine for sending tastes remotely (the user sticks his/her tongue in a device which transforms a signal delivered over the Internet into electrical impulses to the tongue). Coupled with the group's work in developing a machine for sending olfactory signals over a network, the ultimate aim is to build a world repository of gastronomic knowledge, presumably accessible online to users everywhere.

3.2 A Review of Mulsemmedia Research and Applications

Mulsemmedia research is usually inextricably linked to the development of novel and exciting applications. One of the earliest such mulsemmedia VR application is that of Cater (1992) and his team, who developed a virtual reality system to train potential fire-fighters to recognize characteristic smells commonly associated with fires. The problem being solved in this case was to familiarize potential fire-fighters with those smells that are often associated with fires, as it is often thought and argued that it is easier to recognize smells already known by a person. Moreover, in a fire-fighter's profession, being able to detect the presence of such smells could well prove invaluable.

Later on, Dinh et al. (1999) investigated the use of tactile, olfactory and auditory sensory modalities with different levels of visual information on a user's sense of presence and memory of the details of a virtual reality experience. With respect to the olfactory sensory modality, the research study was limited when compared with other sensory modalities considered. Moreover, the single olfactory cue used in the study did not produce any significant effect on the sense of presence, although it did on memory.

One of the benefits of integrating mulsemmedia interfaces in applications is that it can overcome literacy barriers and bring the world of computing closer to categories of people who had hitherto been excluded from it. Jain (2003) was one of the earliest to make this point, when describing the potential of Experiential Computing – computing based on the way humans naturally experience and interact with their environment. Based primarily on video, audio, and taction, he then describes the potential that such interfaces might have in enhancing virtual and augmented reality systems.

In related work, Bodnar, Corbett and Nekrasovski (2004) created a notification system that made use of mulsemmedia data. In their work, they conducted an experimental study to compare the effect of visual, audio or olfactory displays the delivery notifications had on a user's engagement of a cognitive task. Participants were given an arithmetic task to complete and at various intervals two types of notifications were triggered: 1) participants had to immediately stop what they were doing and record some data before returning to the completion of their task, and 2) they could ignore the notification. With this experiment, they found that while olfactory notifications were the least effective in delivering notifications to end users, they had the advantage of producing the least disruptive effect on a user's engagement of a task. It is also worth noting that they encountered most of the problems of using smell output as highlighted earlier by Kaye in their experiment and had participants mostly commenting that some of the smells used were too similar to be distinguishable. Lingering smells in the air also made it difficult to detect the presence of new smells and the lack of experience of working with olfactory data impacted their performance of the assigned task.

Brewster, McGookin and Miller (2006) use explicitly learned odor memories to evaluate the effectiveness of using olfactory data to aid in multimedia content searching, browsing and retrieval in a digital photo library. To conduct this experiment, they developed an olfactory photo browsing and searching tool, which they called Olfoto. The odors are learned by getting participants to complete the explicit odor memory task of associating specific odors with their personal photographs, i.e. smell-based photo tags. Participants were also required to tag the same photographs using text-based tags. The testing phase occurred two weeks later, in which participants were asked to complete three types of exercises. Two were matching exercises that required matching photos with the smell/text tags they had previously associated with them – in one exercise multiple photos were presented with one

smell/text tag and in the other multiple smell/text tags were presented with one photograph. The third exercise involved searching through their digital photo libraries using smell or text tags after being given 3 key features of the photo. Despite the fact that research has shown that odor memories persist longer than word and verbal memories the results showed that performance was lower with the smell-based tags. The lower performance may well be attributed to the fact that possibly odor memories linked to emotions, i.e. those implicitly learned, last longer than those explicitly learned. Thus while in this experimental study participants learned to associate an odor memory with their photographs, the memory was probably not as profound as it would have been if the odor memory had been implicitly learned during the real life moment when the photo was taken. Nonetheless, the findings from their study suggest that odor memories do have the potential to play a role in multimedia content searching.

In related work, the effects of olfaction on information recall in a virtual reality game environment were evaluated by Tortell et al. (2007). In this experimental study, participants engaged in game play in a virtual reality environment. The first phase of the study involved an implicit odor learning period for one group of participants, where subjects had a smell present whilst playing the virtual reality game. The other group of participants in this phase of the experiment had no smell present while they played the game. In the second phase of the experiment, which was an information recall task about the VR environment, participants were again split into two groups. One group performed the task with the same smell that was present during the first phase of the experiment, while the second group performed the information recall task with no smell present. Participants were randomly assigned to groups in the two phases of the experiments, so that participants who completed the first phase of the experiment in the presence of smell did not necessarily get to complete the second phase with the presence of smell and vice versa. Results showed that the subjects who were presented with scent only during the recall phase performed by far the worst, while subjects with scent only during the VR experience performed the best. However, the general findings from the study did show that the introduction of scent in the VR environment had a positive effect on subjects' recollection of the environment.

Multimedia entertainment, such as computer games, is another area that is expected to benefit from the addition of our other sensory cues (thus becoming multimedia games). It is expected that they will heighten the sense of presence and reality and hence impact positively on user experience, e.g. make it a more engaging experience for users. Below, we mention some media entertainment systems that involve the use of olfactory data in one way or another.

Fragra is a Visual-Olfactory virtual reality game that enables players to explore the interactive relationship between olfaction and vision (Mochizuki et al., 2004). The objective of the game is to identify if the visual cues experienced correspond to the olfactory cues at the same time. The game environment has a mysterious tree that bears many kinds of foods. Players can catch these food items by moving their right hand and when they catch one of the items and move it in front of their nose, they smell something which does not necessarily correspond to the food item they are holding. Although they do not report on any detailed evaluation of their implemented game, they do report that in their preliminary experiment, the percentage of questions answered correctly varied according to the combination of visual information and olfactory information and conclude that there is a possibility that some foods' appearance might have stronger information than their scents, and vice versa. A similar interactive computer game, called the "Cooking Game," was created by Nakamoto and his research team at the Tokyo Institute of Technology (Nakamoto et al., 2008).

In earlier related work, Boyd Davis et al. (2006) used olfactory data to create an interactive digital olfactory game. However, the main objective of their experiment, "what should the designer of interactive systems know about olfactory data?" is a question already answered by predecessors in the field. In their work, they developed a suite of digital games in which they use olfactory data, (i.e., three different scents) to engage users in game play. The users' sense of smell is the main skill needed to win the games. The findings from their work further confirm results reported by Kaye about the use of olfactory data.

Morrot et al. (2001) carried out a similar study to investigate the interaction between the vision of colors and odor determination using lexical analysis of wine experts' tasting comments. For the

experiment, they simulate a wine tasting practice, where the wine tasters provide comments on the tasted wines based on the visual, olfactory and gustatory properties of the wines. A previous study (Williams et al., 1984) had actually shown that perception of the olfactory qualities of wines changes depending on whether the color of the wine is visible or hidden from the subjects by using transparent and opaque wine glasses respectively. In the study carried out by Morrot et al., they colored a white wine artificially red and presented it to wine experts to analyze, alongside the uncolored white wine and a red wine. To confirm that the colorant used to artificially color the wine had no influence on the colored wine, a pre-test experiment was carried out to confirm that the white wine and its artificially colored version were perceived as the same when its color was obscured from the tasters. Their results showed that the white wine was perceived as having the odor of a red wine when colored red (all of the wine tasters that participated in the study described the artificially colored wine with terms relating to red wine qualities; the wine's color thus appears to provide significant sensory information, which misleads the subjects' ability to judge flavor; lastly, the mistake is stronger in the presence than in the absence of access to the wine color).

The Research in Augmented & Virtual Environment Systems (RAVES) research group reported a study conducted to investigate the impact of olfaction (concordant and discordant scents) on a user's sense of immersion into a virtual reality environment (Jones et al, 2004). The experimental study involved participants playing a computer game in an immersive virtual experiment. The experimental conditions consisted of a control case where no scents were released while the participant played the game and two experimental cases, one involving concordant scents (e.g., emission of an ocean mist scent as the player passed the ocean and a musty scent when the player was in the fort in the immersive environment) and the other a discordant scent (e.g., smell of maple syrup throughout the game). The results from this study were not statistically significant, however.

It is of little surprise that, because of the relative novelty of the mulsemmedia combinations involved, the studies reviewed so far also explore user acceptance of these new media objects. This is a theme carried forward in more recent research (Ghinea and Ademoye, 2012), which looked at user perception and acceptance of olfactory media combined with the more traditional audio and video.

Kahol et al. (2006) present strategies and algorithms to model context in haptic applications that allow users to explore haptically objects in virtual reality/augmented reality environments. The results from their study show significant improvement in accuracy and efficiency of haptic perception in augmented reality environments when compared to conventional approaches that do not model context in haptic rendering. Indeed, the use of haptics in mulsemmedia VR environments has recently been the subject of other research (as in the work of Apostolopoulos et al., 2012).

In related work, researchers reported on a perceptual study carried out to establish an algorithm to provide high quality inter-media stream synchronization between haptic and audio (voice) media objects in a virtual environment (Ishibashi et al. 2004). Indeed, synchronization seems to be a common theme across mulsemmedia research. Thus, recent work has explored synchronization of olfactory media with audio-visual content (Ghinea and Ademoye, 2010a), while Steinbach et al. (2012) investigated synchronization issues between different modalities and the integration of video and haptics in resource constrained communication networks. Ghinea and Ademoye (2010b) tackled olfaction-enhanced mulsemmedia, by combining computer generated smell with haptic data.

Interactive media and applications have become ubiquitous and compete for attention in our everyday life and work. As discussed by Sarter (2013), this ubiquity has led to an increasing need of effective multimodal interfacing and decisions, including information distribution across different sensory channels to ensure detection, interpretation, and handling of signals. An overview of well-known models of multimodal management was presented by Sarter. In related work, Rob et al. (2013) presented studies of multisensory (audio, tactile, etc.) integration and cross-modal spatial attention to engage more than just a single sense in complex environments. Firstly, multimodal signals were used to reorient spatial attention under the conditions in which unimodal signals may be ineffective. Secondly, multimodal signals are less likely to be masked in noisy environments. And lastly, natural links exist between specific signals and particular behavioral responses. A multimodal system should be designed to minimize any incongruence presented in different sensory modalities that relate to the same event.

We also mention that mulsemmedia has great therapeutic potential. While aromatherapy, music therapy and therapies based on touch all employ primarily one human sense, the creation of multisensory rooms, which give mulsemmedia experiences to individuals with special needs, ranging from learning difficulties to autism, Alzheimer's and dementia, has been reported. Accordingly, the EU Framework Project 5 *MEDIATE* reported research on rooms comprising both visual (e.g. light, color, UV light, projections, illusions), audio (e.g. soothing music), olfactory (i.e. aromatherapy dispensers), and tactile stimuli (i.e., objects with different textures, shapes, vibration) (Gumtau, 2011). Across the Atlantic, and again for therapeutic purposes, *Multisensory Systems* (<http://multisensorysystems.com>) have developed an immersive mulsemmedia system integrating 3D sound, olfaction, vibration and imagery.

Last but not least, mulsemmedia applications were first created in association with the film industry. So it should come as no surprise that the arts and the creative industries continue to experiment mulsemmedia in their content and delivery mechanisms. In so doing, interactive digital experiences are no longer audio-visual creations but mulsemmedia ones. The integration of haptic and olfactory capabilities in many contemporary interactive designs makes the communicative potential of mulsemmedia in terms of sensory, affective, individual and creative expression even more relevant. Thus for instance, Bamboozle theatre (<http://www.bamboozletheatre.co.uk>) and Oily Theatre (<http://www.oilycart.org.uk>) both specialize on multi-sensory performances tailored exclusively for children with autism or complex disabilities. Theatrical mulsemmedia experiences are also for mainstream audiences – Disney's 4D movie experiences featuring tactile and olfactory stimuli on top of the traditional audiovisual presentation have been a staple of audiences for the last 30-40 years. Dynamic Motion Rides (*DyMoRides*) is an Austrian company, who have developed a host of "complex and innovative entertainment attractions," all involving mulsemmedia, for a wide range of entertainment parks worldwide; while the well-known Lowry theatre in Manchester will be staging *Nosferatu* (<http://www.thelowry.com/event/nosferatu>), a mulsemmedia theatrical event in February 2014, no less.

4. MPEG-V: A STANDARD FOR MULSEMEDIA

4.1 Context and Objectives

The initial purpose of the MPEG-V standard was to provide an architecture and associated information representations to enable the interoperability between virtual worlds and the real world. This also explains the name MPEG-V, where "V" stands for virtual world and the standard was entitled "information exchange with virtual worlds", later renamed to "media context and control" to broaden its scope.

The actual architecture of the MPEG-V standard defines interfaces – which are provided in the form of XML- and binary-based representation formats – between digital content providers (incl. virtual worlds) and real-world devices comprising sensors and actuators. These real-world devices may offer various capabilities controlled by appropriate device commands issued by the digital content applications. Alternatively, these commands may be also used to control devices within virtual worlds.

The MPEG-V standard comprises the following parts:

- Part 1: Architecture – describes the general system architecture as well as major interfaces and interoperability points.
- Part 2: Control Information – defines the means to describe the capabilities of (real-world) devices as well as to control them.
- Part 3: Sensory Information – provides the means to describe sensory effects as discussed in the next section.
- Part 4: Virtual World Object Characteristics – provides data representation formats to specify virtual objects that can be exchanged with other virtual worlds.
- Part 5: Data Formats for Interaction Devices – focuses on device interactivity and associated data formats.
- Parts 6 and 7 define common data types and tools needed for the other parts as well as conformance and reference software.

4.2 Sensory Information

The main purpose of MPEG-V Part 3 – Sensory Information – is to enhance both the quality of and user experience of multimedia services by annotating existing multimedia content with additional sensory effects. The main motivation behind this work is that the consumption of multimedia content may stimulate also other human senses – going beyond hearing and seeing – including olfaction, mechanoreception, thermoception, etc. Therefore, multimedia content is annotated providing so-called sensory effects that steer appropriate devices capable of rendering these effects giving the user the sensation of being part of the particular media which results in a worthwhile, informative user experience.

4.2.1 Concept and System Architecture. The concept and system architecture of receiving sensory effects in addition to audio/visual content is depicted in Fig.2. The media and the corresponding sensory effect metadata (SEM) may be obtained from a Digital Versatile Disc (DVD), Blu-ray Disc (BD), or any kind of online service (i.e., download/play or streaming). The media processing engine, which can be deployed on a set-top-box, DVD/BD player, or any other smart device, is responsible for playing the actual media resource and accompanying sensory effects in a synchronized way based on the user's setup in terms of both media and sensory effect rendering. Therefore, the media processing engine may adapt both the media resource and the SEM according to the capabilities of the various rendering devices.

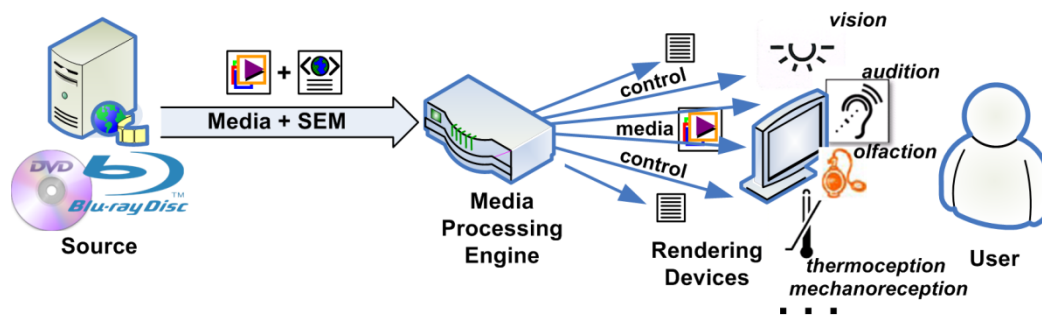


Fig. 2. Concept and System Architecture of Sensory Information

The MPEG-V Part 3 standard deliberately defines only the representation formats without detailing how to create and how to consume multimedia content enriched with sensory effect metadata. This approach enables interoperability among different vendors while supporting a broad range of application domains. Possible means for creating and consuming multimedia with sensory effects including its quality assessment are described in Section 5.

The representation formats defined within MPEG-V Part 3 are now described in the following.

4.2.2 Sensory Effect Description Language (SEDL). The Sensory Effect Description Language (SEDL) is an XML Schema-based language which enables one to describe so-called sensory effects such as light, wind, fog, vibration, etc. that trigger human senses. The actual sensory effects are not part of SEDL but defined within the Sensory Effect Vocabulary (SEV) for extensibility and flexibility allowing each application domain to define its own sensory effects (see Section 4.2.3). A description conforming to SEDL is referred to as Sensory Effect Metadata (SEM) and may be used in any multimedia content (e.g., movies, music, Web sites, games). The SEM can steer sensory devices like fans, vibration chairs, lamps, etc. via an appropriate mediation device to enhance the user experience. That is, in addition to the audio-visual content of, for example, a movie, the user will perceive other effects, giving her/him the sensation of being part of the particular media which should result in a worthwhile, informative user experience.

The current syntax and semantics of SEDL are specified in (Timmerer et al., 2011). However, in this paper we provide an EBNF (Extended Backus–Naur Form)-like overview of SEDL.

```
SEM ::= [autoExtraction] [DescriptionMetadata]
      (Declarations|GroupOfEffects|Effect|ReferenceEffect)+
```

SEM is the root element. It may contain an optional *autoExtraction* and *DescriptionMetadata* attributes followed by a sequence of *Declarations*, *GroupOfEffects*, *Effect*, and *ReferenceEffect* elements. The *autoExtraction* attribute is used to signal whether automatic extraction of a sensory effect from the media resource is preferable. The *DescriptionMetadata* attribute provides information about the SEM itself (e.g., authoring information) and aliases for classification schemes (CS) used throughout the whole description. The MPEG-7 description scheme (Manjunath et al., 2002) is used.

```
Declarations ::= (GroupOfEffects|Effect|Parameter)+
```

The *Declarations* element defines a set of SEDL elements – without instantiating them – for later use in a SEM via an internal reference. In particular, the *Parameter* may be used to define common settings used by several sensory effects similar to variables in programming languages.

A *GroupOfEffects* starts with a *timestamp* that provides information about the point in time when this group of effects should become available for the application. This information can be used for rendering purposes and synchronization with the associated media resource. XML Streaming Instructions as defined in MPEG-21 Digital Item Adaptation (Vetro and Timmerer, 2005) have been adopted for this functionality. Furthermore, a *GroupOfEffects* shall contain at least two *EffectDefinition* for which no timestamps are required as they are provided within the enclosing element. The actual *EffectDefinition* comprises all information pertaining to a single sensory effect.

```
Effect ::= timestamp EffectDefinition
```

An *Effect* is used to describe a single effect with an associated *timestamp*.

```
EffectDefinition ::= [SupplementalInformation] [activate] [duration]
                   [fade-in] [fade-out] [alt] [priority] [intensity] [position]
                   [adaptability] [autoExtraction]
```

An *EffectDefinition* may have a *SupplementalInformation* element for defining a reference region from which the effect information may be extracted in case *autoExtraction* is enabled. Furthermore, several optional attributes are defined which are defined as follows: *activate* describes whether the effect shall be activated; *duration* describes how long the effect shall be activated; *fade-in* and *fade-out* provide means for fading in/out effects respectively; *alt* describes an alternative effect identified by a uniform resource identifier URI (e.g., in case the original effect cannot be processed); *priority* describes the priority of effects with respect to other effects in the same group of effects; *intensity* indicates the strength of the effect in percentage according to a predefined scale/unit (e.g., for wind the Beaufort scale is used); *position* describes the position from where the effect is expected to be received from the user's perspective (i.e., a three-dimensional space is defined in the standard); *adaptability* attributes enable the description of the preferred type of adaptation with a given upper and lower bound; *autoExtraction* with the same semantics as above but only for a certain effect.

4.2.3 *Sensory Effect Vocabulary (SEV)*. The Sensory Effect Vocabulary (SEV) defines a clear set of actual sensory effects to be used with the Sensory Effect Description Language (SEDL) in an extensible and flexible way. That is, it can be easily extended with new effects or by derivation of existing effects thanks to the extensibility feature of XML Schema. Furthermore, the effects are defined in a way to abstract from the authors intention and be independent from the end user's device setting. The sensory effect metadata elements or data types are mapped to commands that control sensory devices based on their capabilities. This mapping is usually provided by the media processing engine and deliberately not defined in this standard, i.e., it is left open for industry competition. It is important to note that there is not necessarily a one-to-one mapping between elements or data types

of the sensory effect metadata and sensory device capabilities. For example, the effect of hot/cold wind may be rendered on a single device with two capabilities, i.e., a heater/air conditioner and a fan/ventilator. Currently, the standard defines the following effects.

Light, colored light, flash light for describing light effects with the intensity in terms of illumination expressed in [lux]. For the color information, a classification scheme (CS) is defined by the standard comprising a comprehensive list of common colors. Furthermore, it is possible to specify the color as RGB. The flash light effect extends the basic light effect by the frequency of the flickering in times per second.

Temperature describes a temperature effect of heating/cooling with respect to the Celsius scale. **Wind** provides a wind effect where it is possible to define its strength with respect to the Beaufort scale. **Vibration** allows one to describe a vibration effect with strength specified using a Richter magnitude scale. For the **water sprayer, scent, and fog** effect the intensity is provided in terms of ml/h.

Finally, the **color correction** effect defines parameters that may be used to adjust the color information in a media resource to the capabilities of end user devices. Furthermore, it is also possible to define a region of interest where the color correction shall be applied in case this desirable (e.g., black/white movies with one additional color such as red).

5. QUALITY OF SERVICE, QUALITY OF EXPERIENCE, AND QUALITY OF SENSORY EXPERIENCE

5.1 Multimedia and Quality of Sensory Experience

New research perspectives on ambient intelligence are presented in Aarts and de Ruyter (2009), which includes also sensory experiences calling for a scientific framework to capture, measure, quantify, judge, and explain the user experience. In a previous paper (de Ruyter and Aarts, 2004) the authors report on the effect additional light effects have on users. User studies showed that light effects are appreciated by users for both audio and visual contents.

In the context of the MPEG-V standardization (Timmerer et al., 2011) some work has been published related to sensory experience that is worth mentioning here. Suk et al. (2009) introduce a new generation of media service called Single Media Multiple Devices (SMMD) which is based on Sensory Effect Metadata (SEM) as defined in MPEG-V. In particular, the SMMD media controller is described that maps sensory effects on appropriate sensory devices for the proper rendering thereof. The main focus of this work is on implementation and engineering. An earlier version puts the controller in the context of Universal Plug and Play (UPnP), thus, focusing also on implementation/engineering aspects (Pyo et al., 2008). Koon et al. (2010) present a framework for 4-D broadcasting based on MPEG-V, that is, the main focus is on delivering additional representation formats in the MPEG-2 Transport Stream (M2TS) and its decoding within the home network environment including the actual service discovery. In this context, Walzl et al. (2013) provide an open-source end-to-end tool chain for creating and consuming multimedia content enriched with sensory effects compliant to MPEG-V based on off-the-shelf infrastructure.

Note that sensory effects are not limited to stationary installations such as in home environments as there is already research to bring sensory effects to mobile devices (Chang and O'Sullivan., 2005). Furthermore, Kim et al. (2010) introduces — among others — new location-based mobile multimedia technology using ubiquitous sensor network-based five senses content. The temporal boundaries within which olfactory data can be used to enhance multimedia applications are investigated in (Ademoye and Ghinea, 2009) concluding that olfaction ahead of multimedia content is more tolerable than olfaction behind content.

Finally, Grega et al. (2008) provide a good overview of the state-of-the-art in QoE evaluation for multimedia services with a focus on subjective evaluation methods which leads us to related work in the area of QoE models. Most of these models focus on a single modality (i.e., audio, image, or video only) or a simple combination of two modalities (i.e., audio and video). For the combination of audio and video content one may employ the basic quality model for multimedia as described in (Hands, 2004). Another approach is known as the IQX hypothesis formulated as an exponential function (Höbßfeld et al., 2008). In (Pereira, 2005) a triple user characterization model for video adaptation and

QoE evaluation is described that introduces at least three quality evaluation dimensions, namely sensorial (e.g., sharpness, brightness), perceptual (e.g., what/where is the content), and emotional (e.g., feeling, sensation) evaluation. Furthermore, it proposes adaptation techniques for the multimedia content and quality metrics associated to each of these layers. The focus is clearly on how an audio/visual resource is perceived, possibly taking into account certain user characteristics (e.g., handicaps) or natural environment conditions (e.g., illumination).

5.2 How to create, consume, and capture QuaSE

In this section, we present a tool chain for creating and consuming media resources annotated with sensory effect including means to capture the Quality of Sensory Experience (QuaSE). This set of tools is one of the first complete end-to-end tool chains offering an easy access from the generation of SEM descriptions till the consumption of audio/video (A/V) content accompanied by SEM descriptions in the context of the World Wide Web or the local playback devices.

Fig. 3 illustrates the whole tool chain starting from the annotation tool (SEVino) on the left side. This tool receives the multimedia content for annotation with sensory effects and outputs the corresponding SEM description. These two assets can then be loaded into the simulator (SESim) located in the center of the figure or delivered via DVD, Blu-Ray, or the Internet. If the content is embedded into a Web site the Web browser plug-in can playback the multimedia content within the Web browser and use the SEM description to steer appropriate devices. If the content is available on other means (e.g., DVD, Blu-Ray) then the stand-alone multimedia player (SEMP) can be used for enhancing the viewing experience. Note that the playback of the Web browser plug-in is performed by the Web browser itself. All tools are freely available under an open-source license and can be downloaded from the Web site of the Sensory Experience Lab (SELab) (<http://selab.itec.aau.at>).

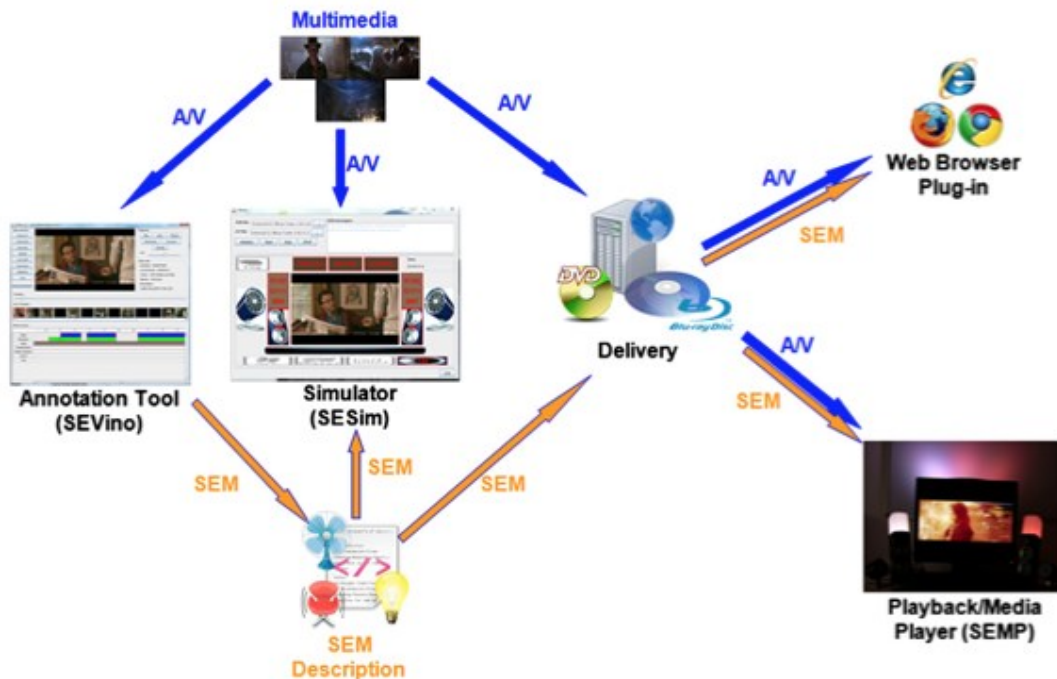


Fig. 3 Overview of end-to-end tool chain enabling to create, consume, and capture QuaSE.

The Sensory Effect Video Annotation (SEVino) tool allows annotating video sequences with various sensory effects (e.g., wind, vibration, light) and generating MPEG-V-compliant SEM descriptions. It is written in Java and for the actual decoding and rendering of the A/V files the Java bindings for VLC¹ are used. Thus, it provides means for embedding the VLC player into a Java application and, thus,

¹ <http://www.videolan.org/vlc/> (last access: March 2014)

enables an application to support a lot of different codecs (e.g., H.264, MPEG-2) and file formats (e.g., MP4, AVI).

The Sensory Effect Simulator (SESim) allows for simulating sensory effects that are contained in SEM descriptions. The Sensory Effect Media Player (SEMP) is a DirectShow-based media player which supports the following devices for rendering sensory effects: the Philips amBX system (with two fans, a wall washer, two light-speakers, a subwoofer, and a wrist rumbler)², the Cyborg Gaming Lights (incl. high-power LEDs)³, and the Vortex Activ device (comprises four slots for providing four different scents)⁴. Note that as the media player uses DirectShow for playback, the media player can handle all formats and codecs which are supported either natively by Windows or via various codec packs.

Finally, the Web browser plugin is based on the AmbientLib which enables arbitrary applications to enrich the user experience with sensory effects. Thus, the library can be seen as an adaptation and processing engine between the virtual description of sensory effects and real devices capable of rendering the described effects. In particular, it provides functionalities to parse SEM descriptions, according to the MPEG-V standard, color calculation of video frames, and enables rendering of sensory effects on a variety of devices. AmbientLib provides an Application Programming Interface (API) and a Driver Interface (DI). The API enables embedding the library within any application and the DI is used for an easy integration of external devices (e.g., those supported also by SEMP) rendering sensory effects. One such application is the Web browser which allows the use of sensory effects with embedded video content on the World Wide Web such as YouTube.

In order to capture the Quality of Experience (QoE) enabled by mulsemmedia, comprising traditional audio-visual content enriched with sensory effects, appropriate subjective quality assessments need to be conducted. Therefore, Walzl et al. (2012) provides a sensory effect dataset and test setups based on the open source tools introduced above. The test setups are aligned with ITU-T's recommendations for subjective quality assessments which provide the basis to study the impact on the QoE when consuming multimedia assets annotated with sensory effects. Timmerer et al. (2012) describes the results of three subjective quality assessments in this domain based on methods defined by ITU-T P.910 and P.911, respectively. The main conclusions from these user studies are that genres such as action, sports, and also documentary benefit from additional sensory effects while the impact on the QoE for genres like commercials and, specifically, news is not that much appreciated. Additionally, media resources with sensory effects may successfully mask visual quality degradations of the actual video content. In the extreme case, the low-quality version of the video enhanced with sensory effects receives higher ratings (on a mean opinion score scale) than the high-quality version of the video with sensory effects. Finally, in (Rainer et al., 2012) the impact on the emotional state is investigated across different sites in Austria and Australia. The results indicate that the intensity of active emotions (e.g., interest, surprise, fun) are increased for video sequences with sensory effects compared to those without sensory effects. The results of the Austrian site also suggest that the intensity of passive emotions (e.g., worry, fear, anger) are decreased for video sequences with sensory effects (compared to those without sensory effects) but with the results from the other sites, it does not yet allow for a general conclusion on whether passive emotions are decreased or increased in their intensity.

Finally, the ultimate goal is to define a utility model which tries to estimate the QoE of multimedia content enhanced with sensory effects based on various influence factors and features. See (Le Callet et al., 2013) for a general definition of QoE. These influence factors and features result from the QoE of the actual multimedia content and the QoE contributions of the individual sensory effects and the combinations thereof. The former can be estimated based on existing models (e.g., such as those referenced in the related work section) whereas the QoE contributions of the sensory effects, both individual and combinations, requires further subjective quality assessments. Therefore, the results of such studies (Walzl et al., 2010; Timmerer et al., 2013) indicate a linear relationship between the

² <http://www.ambx.com/> (last access: March 2014)

³ <http://www.cyborggaming.com/prod/ambx.htm> (last access: March 2014)

⁴ <http://www.daleair.com/vortex-activ> (last access: March 2014)

number of effects and the actual QoE. Thus, the QoE of multimedia content enhanced with sensory effects is referred to as Quality of Sensory Experience (QuaSE) and can be estimated from the QoE of the audio-visual content without sensory effects (QoE_{av}) as depicted as:

$$\text{QuaSE} := QoE_{av} (\delta + \sum w_i b_i)$$

In this utility model, w_i represents the weighting factor for a single sensory effect of type i (i.e., with the given setup as described above, $i \in \{light(l), wind(w), vibration(v)\}$). Additional sensory effect types such as scent may be incorporated easily, e.g., as soon as appropriate devices become available. The variables $b_i \in \{0, 1\}$ depict the binary variables for each effect and are used to indicate whether an effect is present for a given setup. Finally, δ is used for fine-tuning an instantiation of the model.

6. RESEARCH CHALLENGES AND OPEN ISSUES

Mulsemmedia is an emerging and exciting research area that we believe would extract much effort from the related academic and industrial communities. We have pointed out the challenges and possible research work in Appendix A after existing basic technical approaches and computational models are discussed. In this section, we will highlight R&D possibilities for the near future in order to further advance the technology, applications and services, based upon the authors' understanding and project experience in the related fields. Technical advancement is expected to be made in and facilitated by effective algorithm development, substantial database building, meaningful applications and wider user acceptance.

6.1 Mulsemmedia – a solution in search of a killer app?

6.1.1 *Taste – the last frontier?* For computation modeling of the functioning of human senses, as discussed in Section 3.1, most work has been done for audition and vision; significant recent interests have appeared toward olfaction and taction; and gustation is obviously the least investigated topic so far. We expect increasing activities to happen for gustation and the related issues. One challenge that we see is that, since taste buds are located in the mouth, devices that transmit sensations of taste will necessarily be invasive; alternatively, given the close relationship between taste and smell, it would also be interesting to monitor if the solution ultimately adopted will be to use (non-invasive) olfactory inputs to stimulate and engage gustation.

6.1.2 *Attention modeling.* Human attention refers to the cognitive process of selectively concentrating on one aspect of the environment while ignoring other things (Anderson, 2004). As described in section 2, inputs from one sense or different senses compete for human attention. Attention modeling has been formulated as the allocation of processing resources in humans, with a large number of examples in the visual sense (Itti, et al. 1998, Zhang and Lin, 2013) and joint audiovisual senses (Ma, et al. 2005, You, et al. 2007). A comprehensive attention model should evaluate stimuli from all five senses, and this represents a meaningful research challenge for QoE exploration.

6.1.3 *Building databases.* Appropriate databases play important roles in discovering necessary insights for modeling, model parameter determination, and model verification, as evidenced in the related existing visual and audio modeling (Dermot, et al. 2009, Lin and Kuo 2011, Möller, et al. 2011), and cross-database evaluation is essential toward models' generality (Narwaria and Lin 2012, Narwaria, et al. 2012). There have been only a very limited number of databases available for odor (<http://www.odour.org.uk/information.html>, <http://senselab.med.yale.edu/odordb/?db=5>) and touch (<http://brl.ee.washington.edu/HapticsArchive/exp001.html>); more public databases are needed for mulsemmedia (including gustation).

6.1.4 *Mulsemmedia and performing arts/ entertainment.* 4D (and 5D) theatres are a staple attraction of theme parks worldwide and have been imparting 'novel' mulsemmedia experiences to their visitors for some years now. The challenge will be to move such experiences from the theme parks into the mainstream. To some extent this is already happening: vibrating gaming chairs, with integrated

subwoofers (<http://www.4gamers.net/products/ps3/interactive-gaming-chair>), which make users 'feel' the action (and the bass in the audio) are gaining in popularity and becoming more affordable. Nonetheless, in order for mulsemmedia to proliferate in these domains, we need to better understand how audiences react to mulsemmedia effects; this will also enable script authors to effectively integrate them in the respective story lines.

6.1.5 *Mulsemmedia integration, synchronization, and intensities.* Effective integration of mulsemmedia effects requires several questions to be answered: What mulsemmedia combinations work in practice? In what doses/intensities? What synchronization requirements do new media such as olfactory and gustatory media need to satisfy in relation to their counterparts? These are all as-of-yet unanswered questions, which future research needs to target. Once clarified, new – mulsemmedia - authoring tools would need to be written.

6.1.6 *Wearable Mulsemmedia.* The miniaturization of sensors and computing devices alike has led to an increased focus on the potential of wearable technology: recently, both Google (through the Google Glass project - <http://www.google.com/glass/start/>) and Sony (through the SmartWig project <http://www.bbc.co.uk/news/technology-25099262>) have brought to market wearable computing gadgets. If one thinks that individuals already 'wear' perfume and receive vibrating alerts when their smartphones are in silent mode, the potential of wearable devices to transmit mulsemmedia content becomes obvious. Research will need to be done in order to understand how best to integrate such content in wearable devices, and indeed, how best to design such devices so that they can be purveyors of mulsemmedia.

6.1.7 *Mulsemmedia and e-learning.* Mulsemmedia authoring tools would also come in handy for e-learning systems. This, as e-learning systems stand to gain potential benefits from olfaction-enhanced mulsemmedia applications (for instance), as the online learning of certain subject matters, e.g. chemistry, may be further enhanced by the addition of the corresponding smells if it were possible to transmit odors, or more precisely, transmit commands to a smell generating device to mix and emit the required scent over the Internet. Such future work would of necessity need to explore in what contexts and to which extent does mulsemmedia improve communications. In so doing, guidelines about how exactly to use mulsemmedia to achieve a more accurate knowledge transfer would need to be elaborated.

6.1.8 *Mulsemmedia and e-commerce.* The options to feel the texture of a shirt that one wishes to buy, to smell the fragrance one is contemplating of purchasing, of inhaling the aroma, as well as seeing, tasting and experiencing the texture of a gourmet dish before booking a table at the restaurant serving it, all have the potential of moving from the realm of possibilities to that of reality. In so doing, the touch/taste/smell barriers currently characteristic of e-commerce will be overcome.

6.1.9 *User Acceptance and Experience.* We started off this section by highlighting the need for a mulsemmedia killer app. Whilst in the above we have detailed, among others, what we believe to be potentially interesting mulsemmedia developments, we cannot make any predictions for what a killer mulsemmedia app might be. One thing, however, is certain: user acceptance, and more importantly, take-up is essential for any killer app. In order to do this, future work needs to undertake mulsemmedia QoE studies to better understand how mulsemmedia users react to such experiences. Moreover, in so doing, such efforts would also inform the development of objective mulsemmedia QoE metrics.

6.2 Final thought

“Seeing is believing” is an often-quoted idiom. Perhaps not so well known is the fact that the complete idiom, as penned by its author, the 17th century English clergyman, Thomas Fuller, is actually “Seeing is believing, but feeling is the truth.” We subscribe to this statement, but feel that, for mulsemmedia, the idiom is (at least) three sentences too short.

REFERENCES

- Aarts, E. and de Ruyter, B. 2009. New Research Perspectives on Ambient Intelligence, *Journal of Ambient Intelligence and Smart Environments*, 1, 1, 5–14.
- Ademoye, O. and Ghinea, G. 2009. Synchronization of Olfaction-Enhanced Multimedia, *IEEE Transactions on Multimedia*, 11, 3, 561–565.
- Anderson, J. R. 2004. *Cognitive psychology and its implications* (6th ed.). Worth Publishers
- Apostolopoulos, J. G., Chou, P. A., Culbertson, B., Kalker, T., Trott, M. D. and Wee, S., 2012. The Road to Immersive Communication, *Proceedings of the IEEE*, 100, 4, 974–990.
- Ayabe–Kanamura, S., Schicker, I., Laska, M., Hudson, R., Distel, H., Koboyakawa, T., and Saito S., 1998. A Japanese-German cross-cultural study, *Chemical Senses*, 23, 31–38.
- Bodnar, A., Corbett, R. and Nekrasovski, D. 2004. AROMA: Ambient awareness through olfaction in a messaging application: Does olfactory notification make 'scents'?. In *Proceedings Sixth International Conference on Multimodal Interfaces (ICMI'04)*, 183 -- 190.
- Boyd Davis, S., Davies, G., Haddad, R. and Lai, M. 2006. Smell Me: Engaging with an Interactive Olfactory Game. In *Proceedings of the Human Factors and Ergonomics Society 25th Annual Meeting*, 25--40, UK.
- Brewster, S.A., McGookin, D.K. and Miller, C.A. 2006. Olfoto: Designing a smell-based interaction. In *Proceedings CHI 2006: Conference on Human Factors in Computing Systems*, 653 – 662.
- Campbell, D., Jones, E., and Glavin, M., 2009. "Audio quality assessment techniques—A review, and recent developments", *Signal Processing*, 89, 1489–1500.
- Carbon, C.-C. and Jakesch, M. 2013. A Model for Haptic Aesthetic Processing and Its Implications for Design, *Proceedings of the IEEE*, 101, 9, 2123-2133.
- Cater, J.P. 1992. The Nose Have It! Letters to the Editor, *Presence*, 1, 4, 493–494.
- Chang, A. and O'Sullivan, C. 2005. Audio-Haptic Feedback in Mobile Phones, in *Proceedings CHI '05 extended abstracts on Human factors in computing systems, CHI EA '05*, ACM, New York, NY, USA, 2005, 1264–1267.
- Craig, A. D. (2003). Interoception: the sense of the physiological condition of the body. *Current opinion in neurobiology*, 13(4), 500-505.
- Damasio, A. R. 1989. Time-locked multiregional retroactivation: A systems-level proposal for the neural substrates of recall and recognition, *Cognition*, 33, 25–62.
- de Ruyter, B. and Aarts, E., 2004, Ambient Intelligence: Visualizing the Future, in *Proceedings of the working conference on Advanced visual interfaces AVI '04*, ACM Press, New York, NY, USA 203–208.
- DiMaggio, P. 1997. Culture and cognition. *Annual Review Of Sociology*, 23263-287.
- Dinh, H.Q., Walker, N., Hodges, L.F., Song, C. and Kobayashi, A. 1999. Evaluating the importance of multi-sensory input on memory and the sense of presence in virtual environments. In *Proceedings - Virtual Reality Annual International Symposium*, 222--228.
- Fadel, C., & Lemke, C. 2008. *Multimodal learning through media: What the research says*. San Jose, CA: CISCO Systems. Retrieved October, 21, 2010.
- Ghinea, G. and Ademoye, O. 2010a Perceived Synchronization of Olfactory Multimedia , *IEEE Transactions on Systems, Man, and Cybernetics-Part A: Systems and Humans* , 40, 4, 657 – 663.
- Ghinea, G. and Ademoye, O. 2010b. A User Perspective of Olfaction-Enhanced Multimedia. In *Proceedings of the International Conference on Management of Emergent Digital EcoSystems (MEDES '10)*, 277-280, Thailand, Bangkok.
- Ghinea, G. and Ademoye, O., 2011. Olfaction-enhanced multimedia: perspectives and challenges, *Multimedia Tools and Applications*, 55, 3, 601-626.
- Ghinea, G. and Ademoye, O., 2012. The sweet smell of success: Enhancing multimedia applications with olfaction", *ACM Transactions on Multimedia Computing, Communications and Applications* 8, 1, 2.
- Goldstein, E. B. 2013. *Sensation and perception*. Cengage Learning.
- Gray, R., Spence, C., Ho, C. and Tan, H.Z. 2013. Efficient Multimodal Cuing of Spatial Attention *Proceedings of the IEEE*, 101, 9., 2113 – 2121.
- Grega, M., Janowski, L., Leszczuk, M., Romaniak, P., and Papir, Z. 2008 Quality of Experience Evaluation for Multimedia Services - Szacowanie postrzeganej jako ści usług (QoE) komunikacji multimedialnej, *Przegląd Telekomunikacyjny*, 81, 4, 142–153.
- Gumtau, S. 2011. Affordances of touch in multi-sensory embodied interface design. PhD thesis, University of Portsmouth, UK.
- Hands, D. 2004. A Basic Multimedia Quality Model, *IEEE Transactions on Multimedia*, 6, 6, 806–816.
- Heilig, M. L. 1962. Sensorama Simulator, *United States Patent Office (3,050,870)*; Patented August 28, 1962.
- Hinterseer, P. and Steinbach, E. 2006. A psychophysically motivated compression approach for 3D haptic data, *Proc. Int. Symp. Haptic Interfaces Virtual Environ. Teleoperator Syst.*, 35–41.
- Ho, C. and Spence, C., 2005. Olfactory facilitation of dual-task performance, *Neuroscience letters*, 389, 1, 35--40.
- Ishibashi, Y., Kanbara, T., and Tasaka, S., 2004. Inter-stream synchronization between haptic media and voice in collaborative virtual environments. In *Proceedings of the 12th annual ACM international conference on Multimedia*, ACM, New York, NY, USA, 604-611.
- Höfßfeld, T., Hock, D., Tran-Gia, P., Tutschku, K., Fiedler, M. Testing the IQX Hypothesis for Exponential Interdependency between QoS and QoE of Voice Codecs iLBC and G.711, in *Proceedings 18th ITC Specialist Seminar on Quality of Experience*, Karlskrona, Sweden, 2008.
- Itti, L. Koch, C. and Niebur, E. 1998. A model of saliency-based visual attention for rapid scene analysis, *IEEE Trans Patt Anal Mach Intell.*, 20,11, pp. 1254-9.
- ITU-T Rec. P.910, Subjective Video Quality Assessment Methods for Multimedia Applications, April 2008.

- ITU-T Rec. P.911, Subjective Audiovisual Quality Assessment Methods for Multimedia Applications, December 2008.
- Jain, R. 2003. Experiential computing, *Communications of the ACM* 46, 7, 48-55.
- Jayant, N., Johnston, J. and Safranek, R. 1993. Signal compression based on models of human perception, *Proc. IEEE*, 81, 1385–1422.
- Jones, L., Bowers, C.A., Washburn, D., Cortes, A. and Satya, R.V. 2004, The Effect of Olfaction on Immersion into Virtual Environments, in *Human Performance, Situation Awareness and Automation: Issues and Considerations for the 21st Century* Lawrence Erlbaum Associates, 282–285.
- Kahol, K., Tripathi, P., McDaniel, T., Bratton, L. and Panchanathan, S. 2006. Modeling context in haptic perception, rendering, and visualization, *ACM Transactions on Multimedia Computing, Communications and Applications*, 2, 3, 219–240.
- Kammerl, J., Vittorias, I., Nitsch, V., Faerber, B., Steinbach, E., and Hirche, S. 2010. [Perception-based data reduction for haptic force-feedback signals using adaptive deadbands](#), *Presence, Teleoper. Virtual Environ.*, 19, 5, 450–462.
- Kahneman D. 2003. A perspective on judgement and choice. *American Psychologist*. 58, 697-720.
- Kaye, J.N. 2001, *Symbolic Olfactory Display*, Master of Science edn, Massachusetts Institute of Technology, Massachusetts, U.S.A. Available: <http://www.media.mit.edu/~jofish/thesis/>
- Kim, H. Kwon, H.-J., and K.-S. Hong, 2010. Location Awareness-based Intelligent Multi-Agent Technology, *Multimedia Syst.*, 16, (4-5), 275–292.
- Klatzky, R.L., Pawluk, D., and Peer, A. 2013. Haptic Perception of Material Properties and Implications for Applications, *Proceedings of the IEEE*, 101, 9, 2081-2092.
- Le Callet, P., Möller, S. and Perkiš, A. (eds) 2013. Qualinet White Paper on Definitions of Quality of Experience. European Network on Quality of Experience in Multimedia Systems and Services (COST Action IC 1003), Lausanne, Switzerland, Version 1.2.
- Lin, W. 2006. Computational Models for Just-noticeable Difference, Chapter 9 in *Digital Video Image Quality and Perceptual Coding*, eds. H. R. Wu and K. R. Rao, CRC Press.
- Lin, W. and Jay Kuo, C.-C. 2011. Perceptual Visual Quality Metrics: A Survey, *J. of Visual Communication and Image Representation*, 22, 4, 297–312.
- Liu K and Gulliver S. R. 2013. Semiotics in Building space for Working and Living in *Intelligent Building: Design, Management and Operation* ed. Clements-Croome D.
- Lu, Z., Lin, W., Yang, X., Ong, E. and Yao, S. 2005. Modeling Visual Attention's Modulatory Aftereffects on Visual Sensitivity and Quality Evaluation. *IEEE Trans. Image Processing*, 14, 11, 1928 – 1942.
- Ma, Y-F, Hua, X-S , Lu, L. and Zhang, H-J. 2005. A generic framework of user attention model and its application in video summarization. *IEEE Trans. on Multimedia*, 7,5, 907-919.
- Mayer, R. E. 2003. Elements of a science of e-learning. *Journal of Educational Computing Research*, 29(3), 297-313.
- Manjunath, B. S. Salembier, P. and Sikora., T. 2002. Introduction to MPEG-7: Multimedia Content Description Interface, John Wiley and Sons Ltd.
- Marois, R., & Ivanoff, J. 2005. Capacity limits of information processing in the brain. *Trends in cognitive sciences*, 9(6), 296-305.
- Metzinger, T. 1995. Faster than thought. Holism, homogeneity and temporal coding. In T. Metzinger (Ed.), *Conscious experience* 425–461, Paderborn: Schöningh
- Mochizuki, A., Amada, T., Sawa, S., Takeda, T., Motoyashiki, S., Kohyama, K., Imura, M. and Chihara, K. 2004, Fragra: a visual-olfactory VR game. In *Proceedings SIGGRAPH '04: ACM SIGGRAPH 2004 Sketches* ACM Press, New York, NY, USA, pp. 123.
- Möller, S., Chan, W.-Y., Côté, N., Falk, T.H., Raake, A., and Wältermann, M. 2011. Speech Quality Estimation: Models and trends, *IEEE Signal Processing Magazine*, 28, 6, 18–28.
- Morrot, G., Brochet, F. and Dubourdieu, D. 2001, The Color of Odors, *Brain and Language*, 79, 2, 309-320.
- Narwaria, M. and Lin, W. 2012. SVD-Based Quality Metric for Image and Video Using Machine Learning, *IEEE Trans. on Systems, Man, and Cybernetics--Part B*, 42(2), 347 - 364.
- Narwaria, M., Lin, W., McLoughlin, I., Emmanue, S. and Chia, L. T. 2012. Nonintrusive Quality Assessment of Noise Suppressed Speech with Mel-Filtered Energies and Support Vector Regression, *IEEE Trans. on Audio, Speech and Language Processing*, 20(4), 1217 - 1232.
- Nakamoto, T., Otaguro, S., Kinoshita, M., Nagahama, M., Ohinishi, K., and Ishida, T. 2008. Cooking Up an Interactive Olfactory Game Display, *IEEE Computer Graphics and Applications*, 28, 1, 75--78.
- Nothdurft, H.-C. 2000. Saliency from feature contrast: additivity across dimensions. *Vis. Res.*, 40, 10–12, 1183–1201.
- Otaduy, M.A., Garre, C., and Lin, M.C. 2013. Representations and Algorithms for Force-Feedback Display, *Proceedings of the IEEE*, 101, 9, 2068-2080.
- Pereira, F. 2005. A Triple User Characterization Model for Video Adaptation and Quality of Experience Evaluation, in *Proceedings 7th IEEE Workshop on Multimedia Signal Processing*, 1–4.
- Pyo, S., Joo, S., Choi, B., Kim, M., and Kim, J. 2008. A Metadata Schema Design on Representation of Sensory Effect Information for Sensible Media and its Service Framework using UPnP, in *Proceedings 10th International Conference on Advanced Communication Technology, (ICACT 2008)*, 2, 1129 –1134.
- Rainer, B, Walzl, M., Cheng, E., Shujau, M., Timmerer, C., Davis, S., Burnett, I., Ritz, C. and Hellwagner, H. 2012. Investigating the Impact of Sensory Effects on the Quality of Experience and Emotional Response in Web Videos in: I. Burnett, H. Wu (Eds.), *Proceedings of the 4th International Workshop on Quality of Multimedia Experience (QoMEX'12)*, IEEE, Yarra Valley, Australia, 278--283.
- Reinhard, E., Efros, A.A., Kautz, J., and Seidel, H.-P. 2013. On Visual Realism of Synthesized Imagery, *Proceedings of the IEEE*, 101, 9, 1998 -- 2007, 2013.
- Revonsuo, A. 1999. Binding and the phenomenal unity of consciousness. *Consciousness and cognition*, 8, 2, 173-185.
- Richard, G., Sundaram, S., and Narayanan, S. 2013. An Overview on Perceptually Motivated Audio Indexing and Classification,

- Proceedings of the IEEE*, 101, 9, 1939 --1954.
- Rowe, L.A. and Jain, R. 2005, ACM SIGMM retreat report on future directions in multimedia research, *ACM Transactions on Multimedia Computing, Communications, and Applications*, 1, 1, 3--13.
- Rubinstein, J. S., Meyer, D. E., & Evans, J. E. (2001). Executive control of cognitive processes in task switching. *Journal of Experimental Psychology: Human Perception and Performance*, 27(4), 763.
- Sarter, N. 2013. Multimodal Support for Interruption Management: Models, Empirical Findings, and Design Recommendations, *Proceedings of the IEEE*, 101, 9, 2105 – 2112.
- Schiller PH. 1986. The central visual system, *Vision Res.* 26 (9): 1351–1386.
- Seungmoon C. and Kuchenbecker, K.J. 2013. Vibrotactile Display: Perception, Technology, and Applications, *Proceedings of the IEEE*, 101, 9, 2093-2104.
- Smythies, J. R. 1994a. *The walls of Plato's cave*. Aldershot: Avebury.
- Smythies, J. R. 1994b. Requiem for the Identity Theory. *Inquiry*, 37, 311–329.
- Stamper, R.K. 1973 *Information in Business and Administrative Systems*, New York, John Wiley and Sons.
- Steinbach, E., Hirche, S., Ernst, M., Brandi, F., Chaudhari, R., Kammerl, J., and Vittorias, I., 2012. Haptic Communications, *Proceedings of the IEEE*, 100, 4, 937 –956.
- Suk, C. B., Hyun, J. S., and Yong, L. H. , 2009. Sensory Effect Metadata for SMMD Media Service, in *Proceedings of the 2009 Fourth International Conference on Internet and Web Applications and Services*, IEEE Computer Society, Washington, DC, USA, 649–654.
- Timmerer, C, Kim S. K., Ryu, J., and Choi, B. S. 2011. ISO/IEC 23005-3 FDIS Information technology — Media context and control — Part 3: Sensory information.
- Timmerer, C., Walzl, M., Rainer, B., Hellwagner, H. 2012. Assessing the Quality of Sensory Experience for Multimedia Presentations. *Signal Processing: Image Communication* 27, 8, 909--916.
- Tortell, R., Luigi, D.P., Dozois, A., Bouchard, S., Morie, J.F. and Ilan, D. 2007, The effects of scent and game play experience on memory of a virtual environment, *VirtualReality*, 11, 1 , 61–68.
- Vetro, A. and Timmerer, C. 2005. Digital item adaptation: overview of standardization and research activities, *IEEE Transactions on Multimedia*, special issue on MPEG-21, 7, 3, 418 --426.
- Walzl, M., Timmerer, C., Rainer, B., and Hellwagner, H. 2012. Sensory Effect Dataset and Test Setups, in: I. Burnett, H. Wu (Eds.), *Proceedings of the 4th International Workshop on Quality of Multimedia Experience (QoMEX'12)*, IEEE, Yarra Valley, Australia, 115--120.
- Walzl, M., Rainer,B., Timmerer, C., and Hellwagner, H. 2013. An End-to-End tool Chain for Sensory Experience based on MPEG-V., *Signal Processing: Image Communication*, 28, 2, 136--150.
- Williams, A., Langron, S., and Noble, A. 1984. Influence of appearance on the assessment of aroma in Bordeaux wines by trained assessors. *Journal of the Institute of Brewing*, 90, 250–253.
- Wu, H. R., Reibman, A. , Lin, W., Pereira, F., and Hemami S. S. 2013. Perceptual Visual Signal Compression and Transmission, *Proceedings of the IEEE*, 101, 9, 2025 – 2043.
- Yang, X., Lin, W., Lu, Z., Ong, E. and Yao, S. 2005. Just Noticeable Distortion Model and Its Applications in Video Coding. *Signal Processing: Image Communication*, 20, 7, 662-680.
- Yarbus, A. L. 1967. Eye movements during perception of complex objects. In *Eye movements and vision* (pp. 171-211). Springer US.
- Yazdani, A., Kroupi, E., Vesin, J., and Ebrahimi, T., Electroencephalogram alterations during perception of pleasant and unpleasant odors in: I. Burnett, H. Wu (Eds.), *Proceedings of the 4th International Workshop on Quality of Multimedia Experience (QoMEX'12)*, IEEE, Yarra Valley, Australia., 272--277.
- Yoon, K., Choi, B., Lee, E.-S., and Lim, T.-B. 2010. 4-D Broadcasting with MPEG-V, in *Proceedings IEEE International Workshop on Multimedia Signal Processing (MMSP)*, 257 –262.
- You, J., Reiter, U., Hannuksela, M.M., Gabbouj, M., and Perkis, A. 2010. Perceptual-based quality assessment for audio-visual services: a survey, *Signal Processing: Image Communication*, 25, 7, 482–501.
- You, J., Liu, G., Sun, L. and Li, H. 2007. A Multiple Visual Models Based Perceptive Analysis Framework for Multilevel Video Summarization. *IEEE Trans. on Circuits and Systems for Video Technology*, 17,3, 273 – 285.
- Yost, William A. and Nielsen, Donald W., 1985. *Fundamentals of Hearing*, Holt, Rinehart and Winston, New York.
- Zhang, L. M. and Lin, W. 2013. *Modeling Selective Visual Attention: Techniques and Applications*, John Wiley & Sons.

Online Appendix to:

Mulsemmedia: State-of-the- Art, Perspectives and Challenges

GHEORGHITA GHINEA Brunel University
CHRISTIAN TIMMERER Alpen-Adria-Universität
WEISI LIN Nanyang Technological University
STEPHEN R. GULLIVER University of Reading

A. IMPORTANT TECHNICAL ISSUES AND COMPUTATIONAL MODELS IN MULSEMEDIA

This appendix presents more details and discussion for the basic and important technical approaches and computational models for mulsemmedia. We will also try to highlight the related technical challenges and possible future exploration, whenever possible.

A.1 Just noticeable difference (JND) modelling

The just noticeable difference (JND) is the minimum change in the magnitude of a stimulus that can be detected by humans. In a haptic problem, the JND captures certain error tolerance (i.e., the deadband or deadzone) in force and velocity signals below human haptic thresholds, and therefore facilitates effective and efficient data compression (Hinterseer and Steinbach 2006, Kammerl, et al. 2010), error resilience (Steinbach, et al. 2012), rendering (Steinbach, et al. 2012; Seungmoon and Kuchenbecker 2013), interaction and quality evaluation (Kammerl, et al. 2010); for instance, a deadband signal sample needs not to be transmitted to a remote site to achieve computational and bandwidth saving.

In the 1-DoF (degree of freedom) haptic case as in the work by Hinterseer and Steinbach 2006, and Steinbach, et al. 2012, a signal sample $x(t)$ is within the deadband if the following inequality is held:

$$|x(t') - x(t)| \leq k \cdot |x(t')| \quad (\text{A.1})$$

where $x(t')$ is an adjacent sample (usually the previous sample) for $x(t)$, and k is a perceptual threshold parameter simply determined by Weber's law to represent the JND (i.e., the JND is proportional to the signal magnitude) in the aforementioned work.

In real-world haptic systems with multiple DoF, a signal vector $\vec{x} \in R^n$ is used instead of a scalar signal $x(t)$ in (A.1). A multi-DoF signal $\vec{x}(t)$ is within the deadband if:

$$\|\Omega(\vec{x}(t') - \vec{x}(t))\| \leq \|\vec{x}(t')\| \quad (\text{A.2})$$

This is an extension of (A.1) by Steinbach, et al. (2012), and the deadband matrix Ω for n-dimensional signals is:

$$\Omega = \text{diag}\left(\frac{1}{k_1}, \frac{1}{k_2}, \dots, \frac{1}{k_n}\right) \quad (\text{A.3})$$

where k_i is the JND-related threshold corresponding to each element in \vec{x} .

It has been further known that when a human user interacts with an object with a certain velocity $\dot{x}(t)$, his/her force-feedback perception abilities are reduced; that is, the value of k in (A.1) or k_i in (A.3) increases with $\dot{x}(t)$, and therefore now time varying, as in the work by J. Kammerl, et al. (2010):

$$k_v(t) = k_0 + \alpha|\dot{x}(t)| \quad (\text{A.4})$$

where $k_v(t)$ is the velocity-adaptive JND-related threshold, k_0 is the base-line (constant) threshold, and α denotes the rate of change in $k_v(t)$ with respect to $\dot{x}(t)$.

Contrast to audiovisual development for JND (Jayant, et al. 1993, Lin 2006, Lin and Kuo 2011, and Wu, et al. 2013), that for touch sensor and display devices is still in its infancy (e.g., the simple use of Weber's law, as in the existing work mentioned above), while there is lack of similar research in olfaction. Therefore, there is a call for in-depth, comprehensive and systematic investigation for mulsemmedia JND modelling, especially in masking and contrast sensitivity.

A2. Perception of conflicting multisensory information

For a computational model of mulsemmedia, it is inevitable for a discrepancy to occur among different streams of sensory information in space or in time. Some forms of information are of lingering nature (like smell), as opposed to the transitory nature of others (such as video and audio). There has been initial investigation in the related research community, regarding the impact of asynchronisation and need of synchronization of different media.

The olfaction-enhanced multimedia study by Ghinea and Ademoye (2009, 2010a, 2011, 2012) concerns itself with associating computer-generated smell with visual and audio information; the six smell categories used were flowery, foul, fruity, burnt, resinous and spicy, together with the associated videos, as listed in Table A.1 (Ademoye and Ghinea, 2009). Subjective experiments were conducted with more than 40 participants, toward:

- 1) Detectable inter-media skew between olfactory and audiovisual media content;
- 2) Impact of delay on the user-perceived experience.

As shown with the experiments, inter-media skew synchronisation requirements for olfaction and audiovisual content lie between -30 and +20 sec; olfaction ahead of audiovisual content is less noticeable than the reverse case (i.e., olfaction behind audiovisual content). Furthermore, the results revealed that although participants detected the presence of synchronization errors, it did not have a significant impact on the general perceived quality of experience of the olfaction-enhanced multimedia for participants.

Although the human perception system seems to be able to correct for inter-media mismatch so that the discrepancy becomes less noticeable with time, as discussed above, some research indicates that a computational system requiring the user to frequently adapt to novel conflicting situations will have unsatisfactory performance in terms of QoE. Hence, in order to facilitate the coherent perception of an event across different sensory feedbacks, inter-media asynchrony should be systematically minimized in a teleoperation system, for instance, via intelligent statistical multiplexing of audiovisual-haptic signals on the feedback communication channel (Hinterseer and Steinbach, 2006; Kammerl, et al, 2010; Seungmoon and Kuchenbecker, 2013).

Since the perceptual mechanisms behind the conflicting information and synchronization are still largely unknown, obviously more exploration is called for this field, before the findings can be effectively turned into design and implementation advantages for olfaction/haptics-enhanced applications and services.

Table A.1 Associating computer-generated smell with videos (Ademoye and Ghinea, 2009)

SMELL CATEGORY	BURNT	FLOWERY	FOUL	FRUITY	RESINOUS	SPICY
VIDEO DESCRIPTION	Documentary on bush fires in Oklahoma	News broadcast featuring perfume launch	Documentary about rotting fruits	Cookery show on how to make a fruit cocktail	Documentary on Spring allergies& cedar wood	Cookery show on how to make chicken curry
SMELL USED	Burning Wood	Wallflower	Rubbish Acrid	Strawberry	Cedar Wood	Curry

A3. Multisensory integration

Mulsemmedia integration needs to be performed toward the total control and QoE, with $\{s_i\}$ denoting the perceptual effect of noisy multisensory stimuli, where $i=1, 2, \dots, n$, being the sensory index. Assuming noises are independent and Gaussian distributed, a way to integrate the unbiased sensory estimates $\{s_i\}$ is to adaptively weight them (as in the work of haptic communication by Steinbach, et al. (2012)), with w_i to be proportional to the inverse of the variances of respective noise distributions, $\{w_i = 1/\sigma_i^2\}$, i.e.,

$$w_i = \frac{r_i}{\sum_{i=1}^n r_i}, \quad (\text{A.5})$$

and resultant additive integration is:

$$\hat{S} = \sum_{i=1}^n p_i \cdot \hat{s}_i = \sum_{i=1}^n \hat{t}_i \quad (\text{A.6})$$

where p_i is the normalized form of w_i , and $\hat{t}_i = p_i \cdot \hat{s}_i$.

In general, the overlapping effect among $\{s_i\}$ needs to be accounted for, so with extension of the nonlinear additivity model for perception proposed by Nothdurft (2000), Eq. (A.6) becomes

$$\hat{S} = \sum_{i=1}^n \hat{t}_i - \sum_{i,j=1 (i \neq j)}^n \min(c_{ij} \cdot \hat{t}_j, c_{ji} \cdot \hat{t}_i) \quad (\text{A.7})$$

where c_{ij} represents the cross-sensory coupling factors from \hat{s}_j to \hat{s}_i , and is defined in the range of $[0, 1]$ to denote the minimum to the maximum in overlapping, so the second term of the right-hand side of (A.7) accounts for overlapping of multisensory data. Simple examples of using (A.7) can be found in the work of Lu, et al. (2005) and Yang, et al. (2005), for visual signals. Toward immersive environments, further and more convincing research is still a very challenging task for overlapping evaluation and multisensory perceptual fusion, verified with data in big scales.