

Space-variant spatio-temporal filtering of video for gaze visualization and perceptual learning

Citation for published version (APA):

Dorr, M., Jarodzka, H., & Barth, E. (2010). Space-variant spatio-temporal filtering of video for gaze visualization and perceptual learning. In *ETRA '10: Proceedings of the 2010 Symposium on Eye-Tracking Research & Applications* (pp. 307-314). Association for Computing Machinery (ACM).
<https://doi.org/10.1145/1743666.1743737>

DOI:

[10.1145/1743666.1743737](https://doi.org/10.1145/1743666.1743737)

Document status and date:

Published: 01/01/2010

Document Version:

Peer reviewed version

Document license:

CC BY-NC-ND

Please check the document version of this publication:

- A submitted manuscript is the version of the article upon submission and before peer-review. There can be important differences between the submitted version and the official published version of record. People interested in the research are advised to contact the author for the final version of the publication, or visit the DOI to the publisher's website.
- The final author version and the galley proof are versions of the publication after peer review.
- The final published version features the final layout of the paper including the volume, issue and page numbers.

[Link to publication](#)

General rights

Copyright and moral rights for the publications made accessible in the public portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

- Users may download and print one copy of any publication from the public portal for the purpose of private study or research.
- You may not further distribute the material or use it for any profit-making activity or commercial gain
- You may freely distribute the URL identifying the publication in the public portal.

If the publication is distributed under the terms of Article 25fa of the Dutch Copyright Act, indicated by the "Taverne" license above, please follow below link for the End User Agreement:

<https://www.ou.nl/taverne-agreement>

Take down policy

If you believe that this document breaches copyright please contact us at:

pure-support@ou.nl

providing details and we will investigate your claim.

Downloaded from <https://research.ou.nl/> on date: 12 Dec. 2021

Open Universiteit
www.ou.nl



Space-Variant Spatio-Temporal Filtering of Video for Gaze Visualization and Perceptual Learning

Michael Dorr*
Institute for Neuro- and Bioinformatics
University of Lübeck

Halszka Jarodzka†
Knowledge Media Research Center
Tübingen

Erhardt Barth‡
Institute for Neuro- and Bioinformatics
University of Lübeck

Abstract

We introduce an algorithm for space-variant filtering of video based on a spatio-temporal Laplacian pyramid and use this algorithm to render videos in order to visualize pre-recorded eye movements. Spatio-temporal contrast and colour saturation are reduced as a function of distance to the nearest gaze point of regard, i.e. non-fixated, distracting regions are filtered out, whereas fixated image regions remain unchanged. Results of an experiment in which the eye movements of an expert on instructional videos are visualized with this algorithm, so that the gaze of novices is guided to relevant image locations. show that this visualization technique facilitates the novices' perceptual learning.

CR Categories: I.4.9 [Image Processing and Computer Vision]: Applications; I.4.3 [Image Processing and Computer Vision]: Enhancement—Filtering; I.4.10 [Image Processing and Computer Vision]: Image Representation—Multidimensional; K.3.1 [Computers and Education]: Computer Uses in Education—Computer-assisted instruction (CAI)

Keywords: gaze visualization, space-variant filtering, spatio-temporal Laplacian pyramid, perceptual learning

1 Introduction

Humans move their eyes around several times per second to successively sample visual scenes with the high-resolution centre of the retina. The direction of gaze is tightly linked to attention, and what people perceive ultimately depends on where they look [Stone et al. 2003]. Naturally, the ability to record eye movement data led to the need for meaningful visualizations. One-dimensional plots of the horizontal and vertical components of eye position over time have been in use since the very first gaze recording experiments (Delabarre [1898] affixed a small cap on the cornea to transduce eye movements onto a rotating drum, using plaster of Paris as glue). Such plots are useful for detailed quantitative analyses, but not very intuitively interpreted. Other tools supporting interpretation of the data include the visualization of gaze density by means of clustered gaze samples [Heminghaus and Duchowski 2006] or the visualization of other features such as fixation duration [Ramloll et al. 2004].

*e-mail: dorr@inb.uni-luebeck.de

†e-mail: h.jarodzka@iwm-kwmrc.de

‡e-mail: barth@inb.uni-luebeck.de

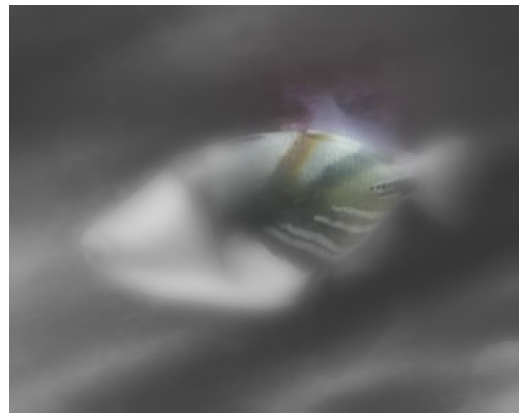


Figure 1: Stillshot from an instructional video on classifying fish locomotion patterns. The eye movements of the expert giving voice-over explanations are visualized by space-variant filtering on a spatio-temporal Laplacian pyramid: spatio-temporal contrast and colour saturation are reduced in unattended areas. This visualization technique aids novices in acquiring the expert's perceptual skills.

Better suited for visual inspection are approaches that use the stimulus and enrich it with eye movement data; in the classical paper of Yarbus [1967], gaze traces overlaid on the original images immediately show the regions that were preferentially looked at by the subjects. Because of the noisy nature of both eye movements and their measurements, there is also an indirect indication of fixation duration (traces are denser in areas of longer fixation). However, such abstract information can also be extracted from the raw data and presented in condensed form: for example, bars of different size are placed in a three-dimensional view of the original stimulus to denote fixation duration in [Lankford 2000]; in a more application-specific manner, Špakov and Rähä [2008] annotate text with abstract information on gaze behaviour for the analysis of translation processes.

Another common method is the use of so-called fixation maps [Velichkovsky et al. 1996; Wooding 2002]. Here, a probability density map is computed by the superposition of Gaussians, each centred at a single fixation (or raw gaze sample), with a subsequent normalization step. Areas that were fixated more often are thus assigned higher probabilities; by varying the width of the underlying Gaussians, it is possible to vary the distance up to which two fixations are considered similar. Based on this probability map, the stimulus images are processed so that for example luminance is gradually reduced in areas that received little attention; so-called heat maps mark regions of interest with transparently overlaid colours. In [Špakov and Miniotas 2007], the authors add “fog” to render visible only the attended parts of the stimulus.

For dynamic stimuli, such as movies, all the above techniques can be applied as well; one straightforward extension from images to image sequences would be to apply the fixation map technique to every video frame individually. Care has to be taken, however, to appropriately filter the gaze input in order to ensure a smooth transition between video frames.

In this paper, we present an algorithm to visualize dynamic gaze density maps by locally modifying spatio-temporal contrast on a spatio-temporal Laplacian pyramid. In regions of low interest, spectral energy is reduced, i.e. edge and motion intensity are dampened, whereas regions of high interest remain as in the original stimulus. Conceptually, this algorithm is related to gaze-contingent displays simulating visual fields based on Gaussian pyramids [Geisler and Perry 2002; Nikolov et al. 2004; Böhme et al. 2006]; in these approaches, however, fine spatial or temporal details are blurred selectively. Instead of blurring, the work presented here leaves details intact but reduces spectral amplitude equally across all frequency bands (note, however, that an individual weighting of separate frequency bands is a trivial extension; also see Section 3.1). Furthermore, while the presented algorithm is based on [Geisler and Perry 2002; Böhme et al. 2006], it cannot be used for gaze-contingent applications where all levels of the underlying pyramid need to be upsampled to full temporal resolution for every video frame. Its purpose is the off-line visualization of pre-recorded gaze patterns.

Pyramid-based rendering as a function of gaze has been shown to have a guiding effect on eye movements [Dorr et al. 2008; Barth et al. 2006]. To further demonstrate the usefulness of our algorithm, we will present some results from a validation experiment in which students received instructional videos either with or without a visualization of the eye movements of an expert watching the same stimulus. Results show that the visualization technique presented here indeed facilitates perceptual learning and improves students' later visual search performance on novel stimuli.

2 Laplacian Pyramid in Space and Space-Time

The so-called Laplacian pyramid serves as an efficient bandpass representation of an image [Burt and Adelson 1983]. In the following section, we will briefly review its application to images and then extend the algorithm to the spatio-temporal domain. We will here use an isotropic pyramid, i.e. all spatial and temporal dimensions are treated equally; this results in a bandpass representation in which e.g. low spatial and low temporal and high spatial and high temporal frequencies are represented together, respectively. For a finer-grained decomposition of the image sequence into spatio-temporal frequency bands, an anisotropic Laplacian pyramid could be used instead. Using such a decomposition, one could also obtain frequency bands of high spatial but low temporal frequencies etc. For a straightforward implementation, one might first create a spatial pyramid for each frame of the input sequence, then decompose each level of that spatial pyramid in time (as in Section 2.2.2, but omitting the spatial up- and downsampling). However, the finer spectral resolution comes at the cost of a significantly increased number of pixels that need to be stored and processed; this increase is on the order of $1.16 \cdot T$ times as many pixels for an anisotropic pyramid with T temporal levels.

2.1 Spatial Domain

The Laplacian pyramid is based on a Gaussian multiresolution pyramid, which stores successively smaller versions of an image; usually, resolution is reduced by a factor of two in each downsam-

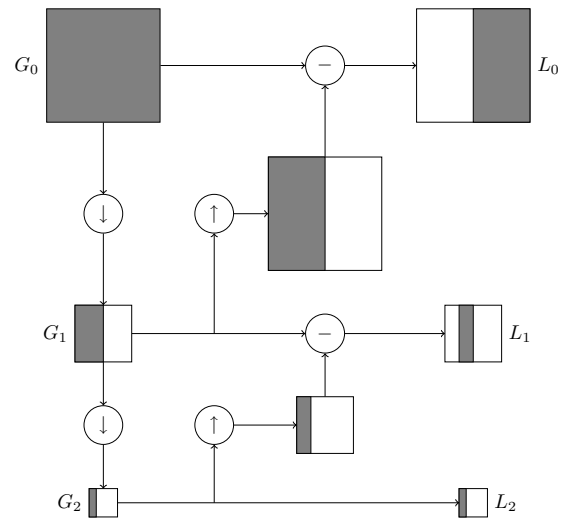


Figure 2: Analysis phase of a Laplacian pyramid in space. Based on the Gaussian pyramid on the left side, which stores successively smaller image versions (with higher-frequency content successively removed), differences of Gaussian pyramid levels are formed to obtain individual frequency bands (right side). To be able to form these differences, lower levels have to be upsampled before subtraction (middle). The gray bars indicate – relative to the original spectrum – what frequency band is stored in each image. The extension into the temporal domain results in lower frame rates for the smaller video versions (not shown).

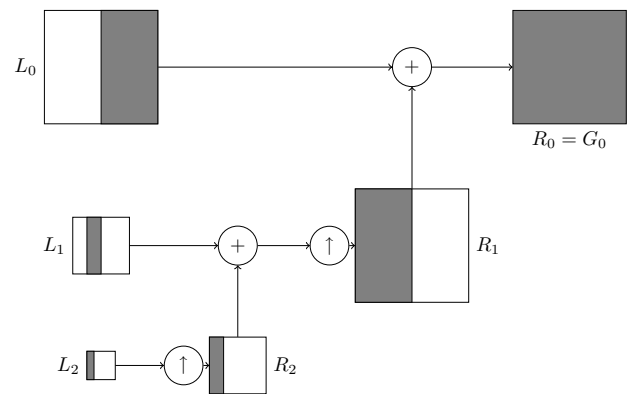


Figure 3: Synthesis phase of a Laplacian pyramid in space. The Laplacian levels are iteratively upsampled to obtain a series of reconstructed images R_N, R_{N-1}, \dots, R_0 with increasing cutoff frequencies. If the L_n remain unchanged, R_0 is an exact reproduction of the original input image G_0 .

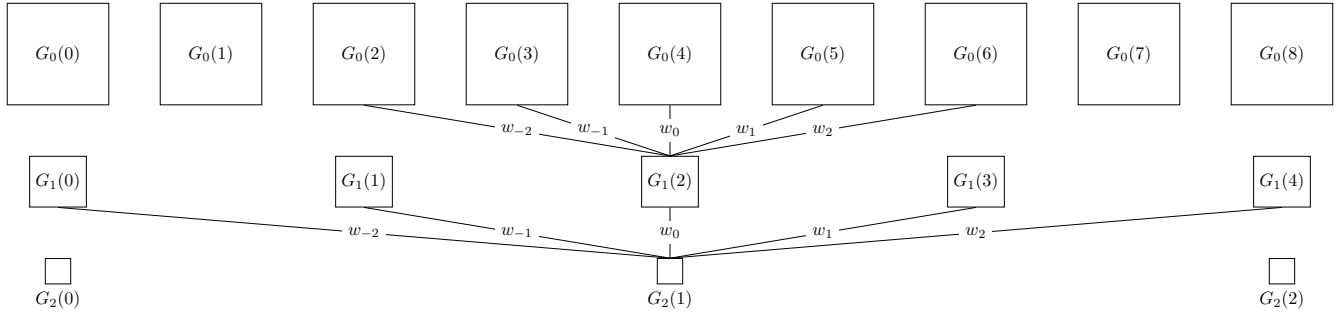


Figure 4: Spatio-temporal Gaussian pyramid with three levels and $c = 2$. Lower levels have reduced resolution both in space and in time. Note that history and lookahead video frames are required because of the temporal filter with symmetric kernel, e.g. computation of $G_2(1)$ depends on $G_1(4)$, which in turn depends on $G_0(6)$.

pling step (for two-dimensional images, the number of pixels is thus reduced by a factor of four). Prior to each downsampling step, the image is appropriately lowpass filtered so that high-frequency content is removed; the lower-resolution downsampling result then fulfills the conditions of the Nyquist theorem and can represent the (filtered) image without aliasing. For a schematic overview, we refer to Figures 2 and 3.

2.2 Spatio-Temporal Domain

The Gaussian pyramid algorithm was first applied to the temporal domain by Uz et al. [1991]. In analogy to dropping every other pixel during the downsampling step, every other frame of an image sequence is discarded to obtain lower pyramid levels (cf. Figure 4). In the temporal domain, however, the problem arises that the number of frames to process is not necessarily known in advance and is potentially infinite; it is therefore not feasible to store the whole image sequence in memory. Nevertheless, video frames from the past and the future need to be accessed during the lowpass filtering step prior to the downsampling operation; thus, care has to be taken which frames to buffer. In the following, we will refer to these frames as history and lookahead frames, respectively.

Both for the subtraction of adjacent Gaussian pyramid levels (to create Laplacian levels) and for the reconstruction step (in which the Laplacian levels are recombined), lower levels first have to be upsampled to match the resolution of the higher level. Following these upsampling steps, the results have to be filtered to interpolate at the inserted pixels and frames; again, history and lookahead video frames are required. We will now describe these operations in more detail and analyse the number of video frames to be buffered.

2.2.1 Notation

The sequence of input images is denoted by $I(t)$; input images have a size of W by H pixels and an arbitrary number of colour channels (individual channels are treated separately). A single pixel at location (x, y) and time t is referred to as $I(t)(x, y)$; in the following, operations on whole images, such as addition, are to be applied pixelwise to all pixels.

The individual levels of a Gaussian multiresolution pyramid with $N + 1$ levels are referred to as $G_i(t)$, $0 \leq i \leq N$. The highest level G_0 is the same as the input sequence; because of the spatio-temporal downsampling, lower levels have fewer pixels and a lower frame rate, so that $G_i(n)$ has a spatial resolution of $W/2^i$ by $H/2^i$ pixels and corresponds to the same point in time as $G_0(2^i n)$. Spatial up- and downsampling operations on an image I are denoted

as $\uparrow [I]$ and $\downarrow [I]$, respectively. For time steps t that are not a multiple of 2^N , not all pyramid levels have a corresponding image $G_i(t/2^i)$; we use C_t to denote the number of levels with valid images at time t (in the implementation, these are the levels that have changed at time t and need to be updated), i.e. C_t is the largest integer with $C_t \leq N$ and $t \bmod 2^{C_t} = 0$. Similar to the Gaussian levels G_i , we refer to the levels of the Laplacian pyramid as $L_i(t)$, $0 \leq i \leq N$ (again, resolution is reduced by a factor of two in all dimensions with increasing i); the intermediate steps during the iterative reconstruction of the original signal are denoted as $R_i(t)$.

The temporal filtering which is required for temporal down- and upsampling introduces a latency (see next sections). The number of lookahead items required on level n is denoted by λ_n for the analysis phase and by Λ_n for the synthesis phase.

2.2.2 Analysis Phase

To compute the Laplacian levels, the Gaussian pyramid has to be created first (see Figure 2). The relationship of different Gaussian levels is shown in Figure 4; lower levels are obtained by lowpass filtering and spatially downsampling higher levels:

$$G_{n+1}(t) = \sum_{i=-c}^c w_i \cdot \downarrow [G_n(2n-i)] \Big/ \sum_{i=-c}^c w_i .$$

We here use a binomial filter kernel $(1, 4, 6, 4, 1)$ with $c = 2$.

The Laplacian levels are then computed as differences of adjacent Gaussian levels (the lowest level L_N is the same as the lowest Gaussian level G_N); before performing the subtraction, the lower level has to be brought back to a matching resolution again by inserting zeros (blank frames) to upsample and lowpass filtering. In practice, the inserted frames can be ignored and their corresponding filter coefficients are set to zero:

$$L_n(t) = G_n(t) - \left[\sum_{i \in P(t)} w_i \cdot G_{n+1} \left(\frac{t-i}{2} \right) \right] \Big/ \sum_{i \in P(t)} w_i ,$$

with $P(t) = \{j = -c, \dots, c \mid (t-j) \bmod 2 = 0\}$ giving the set of valid images on the lower level.

Based on these equations, we can now derive the number of lookahead items required for the generation of the Laplacian. For the upsampling of lower Gaussian levels, we need a lookahead of $\alpha = \lfloor \frac{c+1}{2} \rfloor$ images on each level. Starting on the lowest level G_N , this implies that $2\alpha + c$ images must be available on level G_{N-1}

during the downsampling phase; we can repeatedly follow this argument and obtain $\lambda_n = 2^{N-n} \cdot (\alpha + c) - c$ as the number of required lookahead images for level n .

2.2.3 Synthesis Phase

Turning now to the synthesis phase of the Laplacian pyramid, we note from Figure 3 that the Laplacian levels are successively upsampled and added up to reconstruct the original image; this simply is the inversion of the “upsample-and-subtract” operation during the analysis phase. On the lowest level, $R_N(t) = L_N(t)$; for higher levels, the intermediate reconstructed images are computed as

$$R_n(t) = L_n(t) + \sum_{i \in P(t)} w_i \uparrow \left[R_{n+1} \left(\frac{t-i}{2} \right) \right] \Big/ \sum_{i \in P(t)} w_i .$$

Clearly, a further latency is incurred between the point in time for which bandpass information is available and the reconstructed or filtered image. Similar to the study of the analysis phase in the previous section, we can compute the number Λ_n of required lookahead items on each level by induction. On the lowest level L_N , again $\alpha = \lfloor \frac{c+1}{2} \rfloor$ images are required for the upsampling operation, which corresponds to 2α images on level $N-1$. As can be seen in Figure 5, the result of the upsampling operation is added to the α -th lookahead item on level $N-1$, so that $\Lambda_{N-1} = 3\alpha$. Repeating this computation, we obtain $\Lambda_n = (2^{N+1-n} - 1) \cdot \alpha$; for L_0 , however, no further upsampling is required, so it is possible to reduce the lookahead on the highest level to $\Lambda_0 = (2^{N+1} - 2) \cdot \alpha$.

In practice, we do not need to differentiate explicitly between L and R ; the same buffers can be used for both L and R images. Care has to be taken then not to perform a desired modification of a given Laplacian level on a buffer that already contains information from lower levels as well (i.e. an R image).

2.3 Pseudocode and Implementation

We are now ready to bring together the above observations and put them into pseudocode, see Algorithms 1 and 2. Based on $P(t)$ above, the index function that determines which images are available on lower levels in the following is $P_n(t) = \{j = -c, \dots, c \mid (\frac{t}{2^n} + \alpha - j) \bmod 2 = 0\}$. In the synthesis phase, the image offset α at which the recombination L and R takes place can be set to 0 on the highest level; we therefore use $\alpha_n = \alpha$ for $n > 0$, $\alpha_0 = 0$.

t Time step to update the pyramid for
Input: G_0, \dots, G_N Levels of the Gaussian pyramid
 L_0, \dots, L_N Levels of the Laplacian pyramid

// Gaussian pyramid creation
 $C_t = \max(\{\gamma \in \mathbb{N} \mid 0 \leq \gamma \leq N, t \bmod 2^\gamma = 0\})$
 $G_0(t + \lambda_0) = I(t + \lambda_0)$
for $n = 1, \dots, C_t$ **do**

$$G_n \left(\frac{t}{2^n} + \lambda_n \right) =$$

$$\sum_{i=-c}^c w_i \cdot \downarrow \left[G_{n-1} \left(\frac{t}{2^{n-1}} + 2\lambda_n - i \right) \right] \Big/ \sum_{i=-c}^c w_i$$

end

// Laplacian pyramid creation

for $n = 0, \dots, C_t$ **do**

if $n = N$ **then**

$$L_N \left(\frac{t}{2^N} + \Lambda_N \right) = G_N \left(\frac{t}{2^N} + \Lambda_N \right)$$

else

$$L_n \left(\frac{t}{2^n} + \Lambda_n \right) = G_n \left(\frac{t}{2^n} + \Lambda_n \right) -$$

$$\uparrow \left[\sum_{i \in P_n(t)} w_i \cdot G_{n+1} \left(\frac{t}{2^{n+1}} + \frac{\Lambda_n - i}{2} \right) \right] \Big/ \sum_{i \in P_n(t)} w_i$$

end

end

Algorithm 1: Pseudocode for one time step of the pyramid analysis phase.

From the pseudocode, a buffering scheme for the implementation directly follows. First, images from the Gaussian pyramid have to be stored; each level n needs at least λ_n lookahead images, one current image, and α_n history. Trading memory requirements for computational costs, it is also possible to keep all images of the Gaussian pyramid in memory twice, once in the “correct” size and once in the downsampled version; for each frame of the input video, only one downsampling operation has to be executed then. In analogy to the Gaussian levels, both the Laplacian and the (partially) reconstructed levels L and R can be held together in one buffer per level n with Λ_n lookahead, one current image, and α_n history.

In practice, input images are fed into lookahead position λ_0 of buffer G_0 , and images are shifted towards the “current” position by one position for every new video frame. This means that λ_0 many time steps after video frame I has been added, the corresponding Gaussian images G_0 to G_N are stored in the “current” positions of the Gaussian buffers. The resulting differences L_0 to L_N then are stored at the lookahead positions Λ_0 to Λ_N of the Laplacian buffers, respectively; here, different frequency bands can be accessed both for analysis and modification. Only Λ_0 time steps later does the input image I re-appear after pyramid synthesis; overall, this leads to a pyramid latency between input and output of $\lambda_0 + \Lambda_0$ time steps.

The necessary buffering and the handling of lookahead frames could be reduced and simplified if non-causal filters were used; a further possibility to efficiently filter in time without lookahead is to use temporally recursive filters. However, any non-causality in the filters will introduce phase shifts. Particularly in the case of space-variant filtering (see below), this would produce image artefacts (such as a pedestrian with disconnected – fast – legs and – relatively slow – upper body).

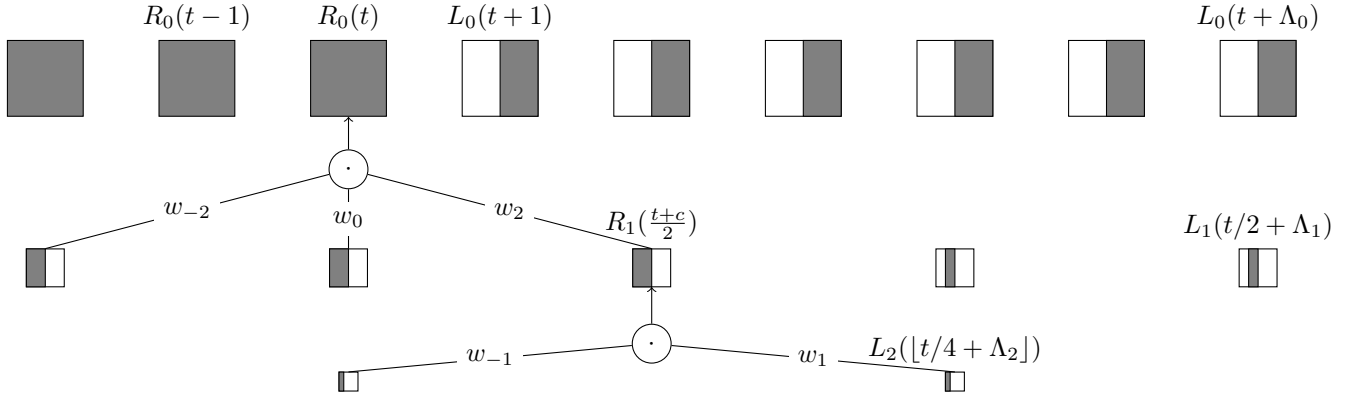


Figure 5: Synthesis step of spatio-temporal Laplacian pyramid. Here shown are both the Laplacian levels L_i and the (partially) reconstructed images R_i , which are based on lower levels with indices $\geq i$; in practice, the same buffers can be used for both R and L . For example, to compute the reconstruction $R_0(t)$ of the original image, we have to add $L_0(t)$ to a spatio-temporally upsampled version of R_1 . The second level L_1 is combined with the upsampling result of L_2 in $R_1(\lfloor \frac{t+c}{2} \rfloor) = R_1(\frac{t}{2} + \alpha_1)$ (see pseudocode). In this schematic overview, new frames are added on the right side and shifted leftwards with time.

t Time step to update the pyramid for
Input: G_0, \dots, G_N Levels of the Gaussian pyramid
 L_0, \dots, L_N Levels of the Laplacian pyramid
 $C_t = \max(\{\gamma \in \mathbb{N} \mid 0 \leq \gamma \leq N, t \bmod 2^\gamma = 0\})$
for $n = C_t - 1, \dots, 0$ **do**
 if $n = N - 1$ **then**
 $R_N(\frac{t}{2^N} + \alpha_N) = L_N(\frac{t}{2^N} + \alpha_N)$

end

$$R_n\left(\frac{t}{2^n} + \alpha_n\right) = L_n\left(\frac{t}{2^n} + \alpha_n\right) + \left[\sum_{i \in P_n(t)} w_i \cdot R_{n+1}\left(\frac{t}{2^{n+1}} + \frac{\alpha_n - i}{2}\right) \right]$$

end

Algorithm 2: Pseudocode for one time step of the pyramid synthesis phase.

3 Gaze Visualization

In the previous section, we described the analysis and synthesis phase of a spatio-temporal Laplacian pyramid. However, the result of the synthesis phase is a mere reconstruction of the original image sequence; we want to filter the image sequence based on a list of gaze positions instead.

3.1 Space-Variant Pyramid Synthesis

We now introduce the concept of a *resolution map* that indicates how spectral energy should be modified in each frequency band at each pixel of the output image sequence. We denote the resolution map for level n at time t with $W_n(t)$; the W_n have the same spatial resolution as the corresponding L_n , i.e. $W/2^n$ by $H/2^n$ pixels.

To bandpass-filter the image sequence, the Laplacian levels L_n are simply multiplied pixel-wise with the W_n prior to the recombination into R_n .

Based on the pseudocode, we can see that resolution maps for dif-

ferent points in time are applied to the different levels in each synthesis step of the pyramid; this follows from the iterative recombination of L into the reconstructed levels. In practice, a more straightforward solution is to apply resolution maps corresponding to one time t to the farthest lookahead item Λ_n of each level L (i.e. right after subtraction of adjacent Gaussian levels).

As noted before, in the following validation experiment we will use the same resolution map for all levels (for computational efficiency, however, resolution maps for lower levels can be stored with fewer pixels). In principle, this means that a similar effect could be achieved by computing the mean pixel intensity of the whole image sequence and then, depending on gaze position, smoothly blending between this mean value and each video pixel. However, for practical reasons, the lowest level of the pyramid does not represent the “true” DC (the mean of the image sequence), but merely a very strongly lowpass-filtered video version; this means that some coarse spatio-temporal structure remains even in regions where all contrast in higher levels is removed by setting the resolution map to zero. The temporal multi-resolution character of the pyramid also adds smoothness to changes in the resolution maps over time; because temporal levels are updated at varying rates, such changes are introduced gradually. Finally, by using different resolution maps for each level, it is trivially possible to highlight certain frequency bands, which is impossible based on a computation of the mean alone.

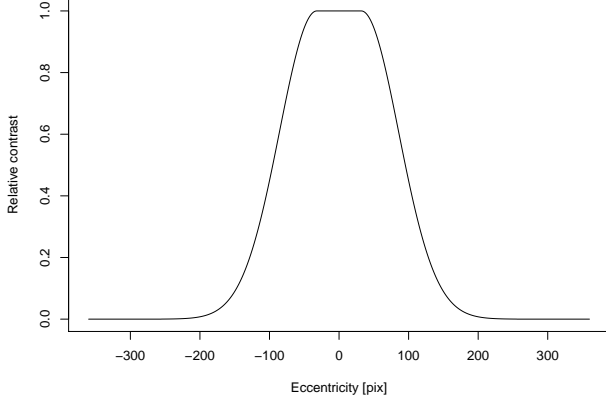


Figure 6: Eccentricity-dependent resolution map: at centre of fixation, spectral energy remains the same; energy is increasingly reduced with increasing distance from gaze position.

Input:

t Time step to update the pyramid for
 G_0, \dots, G_N Levels of the Gaussian pyramid
 L_0, \dots, L_N Levels of the Laplacian pyramid
 $W_n(t)$ Resolution map for frequency band n at time t

$C_t = \max(\{\gamma \in \mathbb{N} \mid 0 \leq \gamma \leq N, t \bmod 2^\gamma = 0\})$

for $n = C_t - 1, \dots, 0$ **do**

if $n = N - 1$ **then**

$$R_N\left(\frac{t}{2^N} + \alpha_N\right) = L_N\left(\frac{t}{2^N} + \alpha_N\right)$$

end

$$R_n\left(\frac{t}{2^n} + \alpha_n\right) = W\left(x, y, \frac{t}{2^n} + \alpha_n\right) \cdot L_n\left(x, y, \frac{t}{2^n} + \alpha_n\right) +$$

$$\left[\sum_{i \in P_n(t)} w_i \cdot R_{n+1}\left(\frac{t}{2^{n+1}} + \frac{\alpha_n - i}{2}\right) \right]$$

end

Algorithm 3: Pseudocode for one time step of the space-variant pyramid synthesis phase.

4 Attentional Guidance

Pyramid-based rendering of video as a function of gaze has been shown to have a guiding effect on eye movements. For example, the introduction of peripheral temporal blur on a gaze-contingent display reduces the number of large-amplitude saccades [Barth et al. 2006], even though the visibility of such blur is low [Dorr et al. 2005]. Using a real-time gaze-contingent version of a spatial Laplacian pyramid (which is computationally much cheaper than the spatio-temporal pyramid presented here), locally reducing (spatial) spectral energy at likely fixation points also changes eye movement characteristics [Dorr et al. 2008].

In the following, we will therefore briefly summarize how our gaze visualization algorithm can be applied in a learning task to guide the student’s gaze. For further details of this experiment, we refer to [Jarodzka et al. 2009c].

4.1 Perceptual Learning

In many problem domains, experts develop efficient eye movement strategies because the underlying problem requires substantial visual search. Examples include the analysis of radiograms [Lesgold et al. 1988], driving [Underwood et al. 2003], and the classification of fish locomotion [Jarodzka et al. 2009a]. In order to aid novices in acquiring the efficient eye movement strategies of an expert, it is possible to use cueing to guide their attention towards relevant stimulus locations; however, it often remains unclear where and how to cue the user. Van Gog et al. [2009] guided attention during problem-solving tasks by directly displaying the eye movements of an expert made during performing the same task on modeling examples, but found that the attentional guidance actually decreased novices’ subsequent test performance instead of facilitating the learning process. One possible explanation of this effect could be that the chosen method of guidance (a red dot at the experts’ gaze position that grew in size with fixation duration) was not optimal because the gaze marker covered exactly the visual features it was supposed to highlight, and its dynamical nature might have distracted the observers. To avoid this problem, we here use the space-variant filtering algorithm presented in the previous sections to render instructional videos so that the viewer’s attention is guided to those areas that were attended by the expert. However, instead of altering these areas, we decrease spatio-temporal contrast (i.e. edge and motion intensity) elsewhere, in order to increase the relative visual saliency of the problem-relevant areas without covering them or introducing artefacts.

4.2 Stimulus Material and Experimental Setup

Twenty-six videos of different fish species with a duration of 4 s each were recorded, depicting different locomotion patterns. They had a spatial resolution of 720 by 576 pixels and a frame rate of 25 frames per second. Four of these videos were shown in a continuous loop to an expert on fish locomotion (a professor of marine zoology) and his eye movements were collected using a Tobii 1750 remote eye tracker running at 50 Hz. Simultaneously, a spoken didactical explanation of the locomotion pattern (i.e. how different body parts moved) was recorded. These four videos were shown to 72 subjects (university students without prior task experience) in a training phase either as-is or with attentional guidance (either a simple yellow square at gaze position or the pyramid-based contrast reduction, see below); the remaining 22 videos served as test videos and were subsequently shown to the subjects, who had to name and describe the locomotion pattern displayed in each test video. In order to measure the similarity of subjects’ eye movements to those of the expert, gaze data was also recorded during training and test phases.

4.3 Gaze Filtering

Functionally, a sequence of eye movements consists of a series of fixations, where eye position remains constant, and saccades, during which eye position changes rapidly (smooth pursuit movements here can be understood as fixations where gaze position remains constant on a moving object). In practice, however, the eye position as measured by the eye tracker hardly ever stays constant from one sample to the next; the fixational instability of the oculomotor system, minor head movements, and noise in the camera system of the eye tracker all contribute to the effect that the measured eye position exhibits a substantial jitter. If this jitter were to be replayed to the novice, such constant erratic motion might distract the observer from the very scene that gaze guidance is supposed to highlight. In order to reduce the jitter, raw gaze data was filtered with a temporal Gaussian lowpass filter with a support of 200 ms and a standard

deviation of 42 ms.

4.4 Space-Variant Filtering and Colour Removal

A Laplacian pyramid with five levels was used; resolution maps were created in such a way that the original image sequence was reconstructed faithfully in the fixated area (the weight of all levels during pyramid synthesis was set to 1.0) and spatio-temporal changes were diminished (all level weights set to 0.0) in those areas that the expert had only seen peripherally. On the highest level, the first zone was defined by a radius of 32 pixels around gaze position and weights were set to 0.0 outside a radius of 256 pixels; these radii approximately corresponded to 1.15 and 9.2 degrees of visual angle, respectively. In parafoveal vision, weights were gradually decreased from 1.0 to 0.0 for a smooth transition, following a Gaussian falloff with a standard deviation of 40 pixels (see Figure 6). Furthermore, these maps were produced not only by placing a mask at the current gaze position in each video frame; instead, masks for all gaze positions of the preceding and following 300 ms were superimposed and the resolution map was then normalized to a maximum of 1.0. During periods of fixation, this superposition had little or no effect; during saccades, however, this procedure elongated the radially symmetric resolution map along the direction of the saccade. Thus, the observer was able to follow the expert's saccades and unpredictable large displacements of the unmodified area were prevented. Finally, colour saturation was also removed from non-attended areas similar to the reduction of spectral energy; here, complete removal of colour started outside a radius of 384 pixels around gaze, and the Gaussian falloff in the transition area had a standard deviation of 67 pixels. Note that these parameters were determined rather informally to find a reasonable trade-off between a focus that would be too restricted (if the focus were only a few pixels wide, discrimination of relevant features would be impossible) and a wide focus that would be without effect (if the unmodified area encompassed the whole stimulus). As such, these parameters are likely to be specific to the stimulus material used here; a systematic analysis might yield better results. For a thorough investigation of the visibility of peripherally removed colour saturation using a gaze-contingent display, we refer to Duchowski et al. [2009].

For an example frame, see Figure 1; a demo video is available online at <http://www.gazecom.eu/demo-material>.

4.5 Results

TODO elaborate ... Gaze marker facilitates perceptual learning: subjects look at relevant movie regions more and take less time to find relevant locations after stimulus onset [Jarodzka et al. 2009b]. The gaze visualization technique presented here does not cover these relevant locations; subjects' visual search performance is improved even beyond that obtained with the simple gaze marker (50% faster onset).

For a more in-depth analysis, we refer to [Jarodzka et al. 2009c].

5 Conclusion

We have presented a novel algorithm to perform space-variant filtering of a movie based on a spatio-temporal Laplacian pyramid. One application is the visualization of eye movements on videos; spatio-temporal contrast is modified as a function of gaze density, i.e. spectral energy is reduced in regions of low interest. In a validation experiment, subjects watched instructional videos on fish locomotion either with or without visualization of the eye movements of an expert. We were able to show that on novel test stimuli, subjects who had received such information performed better

than subjects who had not seen the expert's eye movements during training, and that the gaze visualization technique presented here facilitated learning better than a simple gaze display (yellow gaze marker). In principle, any visualization technique that reduces the relative visibility of those regions not attended by the expert might have a similar effect; our choice for this particular technique was motivated by our work on eye movement prediction [Dorr et al. 2008; Vig et al. 2009], which shows that spectral energy is a good predictor for eye movements. Ultimately, we intend to use similar techniques in a gaze-contingent fashion in order to guide the gaze of an observer [Barth et al. 2006].

Acknowledgements

Our research has received funding from the European Commission within the project GazeCom (contract no. IST-C-033816) of the 6th Framework Programme. All views expressed herein are those of the authors alone; the European Community is not liable for any use made of the information. Halszka Jarodzka was supported by the Leibniz-Gemeinschaft within the project "Resource-adaptive design of visualizations for supporting the comprehension of complex dynamics in the Natural Sciences".

References

- BARTH, E., DORR, M., BÖHME, M., GEGENFURTNER, K. R., AND MARTINETZ, T. 2006. Guiding the mind's eye: improving communication and vision by external control of the scanpath. In *Human Vision and Electronic Imaging*, B. E. Rogowitz, T. N. Pappas, and S. J. Daly, Eds., vol. 6057 of *Proc. SPIE*. Invited contribution for a special session on Eye Movements, Visual Search, and Attention: a Tribute to Larry Stark.
- BÖHME, M., DORR, M., MARTINETZ, T., AND BARTH, E. 2006. Gaze-contingent temporal filtering of video. In *Proceedings of Eye Tracking Research & Applications (ETRA)*, 109–115.
- BURT, P. J., AND ADELSON, E. H. 1983. The Laplacian pyramid as a compact image code. *IEEE Transactions on Communications* 31, 4, 532–540.
- DELABARRE, E. B. 1898. A method of recording eye-movements. *American Journal of Psychology* 9, 4, 572–574.
- DORR, M., BÖHME, M., MARTINETZ, T., AND BARTH, E. 2005. Visibility of temporal blur on a gaze-contingent display. In *APGV 2005 ACM SIGGRAPH Symposium on Applied Perception in Graphics and Visualization*, 33–36.
- DORR, M., VIG, E., GEGENFURTNER, K. R., MARTINETZ, T., AND BARTH, E. 2008. Eye movement modelling and gaze guidance. In *Fourth International Workshop on Human-Computer Conversation*.
- DUCHOWSKI, A. T., BATE, D., STRINGFELLOW, P., THAKUR, K., MELLO, B. J., AND GRAMOPADHYE, A. K. 2009. On spatiochromatic visual sensitivity and peripheral color LOD management. *ACM Transactions on Applied Perception* 6, 2, 1–18.
- GEISLER, W. S., AND PERRY, J. S. 2002. Real-time simulation of arbitrary visual fields. In *ETRA '02: Proceedings of the 2002 symposium on Eye tracking research & applications*. ACM, New York, NY, USA, 83–87.
- HEMINGHOUS, J., AND DUCHOWSKI, A. T. 2006. iComp: a tool for scanpath visualization and comparison. In *APGV '06: Proceedings of the 3rd symposium on Applied perception in graphics and visualization*, ACM, New York, NY, USA, 152–152.

- JARODZKA, H., SCHEITER, K., GERJETS, P., AND VAN GOG, T. 2009. In the eyes of the beholder: How experts and novices interpret dynamic stimuli. *Journal of Learning and Instruction*. (to appear).
- JARODZKA, H., SCHEITER, K., GERJETS, P., VAN GOG, T., AND DORR, M. 2009. How to convey perceptual skills by displaying experts' gaze data. In *Proceedings of the 31st Annual Conference of the Cognitive Science Society*, Austin, TX: Cognitive Science Society, N. A. Taatgen and H. van Rijn, Eds., 2920–2925.
- JARODZKA, H., VAN GOG, T., DORR, M., SCHEITER, K., AND GERJETS, P. 2009. Guiding attention guides thought, but what about learning? Eye movements in modeling examples. (submitted).
- LANKFORD, C. 2000. Gazetracker: software designed to facilitate eye movement analysis. In *ETRA '00: Proceedings of the 2000 symposium on Eye tracking research & applications*, ACM, New York, NY, USA, 51–55.
- LESGOLD, A., RUBINSON, H., FELTOVICH, P., GLASER, R., KLOPFER, D., AND WANG, Y. 1988. Expertise in a complex skill: diagnosing X-ray pictures. In *The nature of expertise*, M. T. H. Chi, R. Glaser, and M. Farr, Eds. Hillsdale, NJ: Erlbaum, 311–342.
- NIKOLOV, S. G., NEWMAN, T. D., BULL, D. R., CANAGARAJAH, C. N., JONES, M. G., AND GILCHRIST, I. D. 2004. Gaze-contingent display using texture mapping and OpenGL: system and applications. In *Eye Tracking Research & Applications (ETRA)*, 11–18.
- RAMLOLL, R., TREPAGNIER, C., SEBRECHTS, M., AND BEEDASY, J. 2004. Gaze data visualization tools: Opportunities and challenges. In *IV '04: Proceedings of the Information Visualisation, Eighth International Conference*, IEEE Computer Society, Washington, DC, USA, 173–180.
- STONE, L. S., MILES, F. A., AND BANKS, M. S. 2003. Linking eye movements and perception. *Journal of Vision* 3, 11 (11), i–iii.
- UNDERWOOD, G., CHAPMAN, P., BROCKLEHURST, N., UNDERWOOD, J., AND CRUNDALL, D. 2003. Visual attention while driving: Sequences of eye fixations made by experienced and novice drivers. *Ergonomics* 46, 629–646.
- UZ, K. M., VETTERLI, M., AND LEGALL, D. J. 1991. Interpolative multiresolution coding of advanced television with compatible subchannels. *IEEE Transactions on Circuits and Systems for Video Technology* 1, 1, 86–99.
- VAN GOG, T., JARODZKA, H., SCHEITER, K., GERJETS, P., AND PAAS, F. 2009. Effects of attention guidance during example study by showing students the models' eye movements. *Computers in Human Behavior* 25, 785–791.
- VELICHKOVSKY, B., POMPLUN, M., AND RIESER, J. 1996. Attention and communication: Eye-movement-based research paradigms. In *Visual Attention and Cognition*, W. H. Zangemeister, H. S. Stiehl, and C. Freksa, Eds. Amsterdam, Netherlands: Elsevier Science, 125–54.
- VIG, E., DORR, M., AND BARTH, E. 2009. Efficient visual coding and the predictability of eye movements on natural movies. *Spatial Vision* 22, 5, 397–408.
- ŠPAKOV, O., AND RÄIHÄ, K.-J. 2008. KiEV: A tool for visualization of reading and writing processes in translation of text. In *ETRA '08: Proceedings of the 2008 symposium on Eye tracking research & applications*, ACM, New York, NY, USA, 107–110.
- WOODING, D. S. 2002. Fixation maps: Quantifying eye-movement traces. In *ETRA '02: Proceedings of the 2002 symposium on Eye tracking research & applications*, ACM, New York, NY, USA, 31–36.
- YARBUS, A. L. 1967. *Eye Movements and Vision*. Plenum Press, New York.