



Office of Graduate Studies

Dissertation / Thesis Approval Form

This form is for use by all doctoral and master's students with a dissertation/thesis requirement. Please print clearly as the library will bind a copy of this form with each copy of the dissertation/thesis. All doctoral dissertations must conform to university format requirements, which is the responsibility of the student and supervising professor. Students should obtain a copy of the Thesis Manual located on the library website.

Dissertation/Thesis Title: New Methods for Detecting Frame Deletion in Modern Video

Author: Hunter Kippen

This dissertation/thesis is hereby accepted and approved.

Signatures:

Examining Committee

Chair

Members

Academic Advisor

Department Head

New Methods for Detecting Frame Deletion in Modern Video

A Thesis

Submitted to the Faculty

of

Drexel University

by

Hunter Kippen

in partial fulfillment of the

requirements for the degree

of

Master of Science

May 2019



© Copyright 2019
Hunter Kippen. All Rights Reserved.

This work is licensed under the terms of the Creative Commons Attribution-ShareAlike
4.0 International license. The license is available at
<http://creativecommons.org/licenses/by-sa/4.0/>.

Dedications

Hello, here is some text without a meaning. This text should show what a printed text will look like at this place. If you read this text, you will get no information. Really? Is there no information? Is there a difference between this text and some nonsense like “Huardest gefburn”? Kjift – not at all! A blind text like this gives you information about the selected font, how the letters are written and an impression of the look. This text should contain all letters of the alphabet and it should be written in of the original language. There is no need for special content, but the length of words should match the language.

This is the second paragraph. Hello, here is some text without a meaning. This text should show what a printed text will look like at this place. If you read this text, you will get no information. Really? Is there no information? Is there a difference between this text and some nonsense like “Huardest gefburn”? Kjift – not at all! A blind text like this gives you information about the selected font, how the letters are written and an impression of the look. This text should contain all letters of the alphabet and it should be written in of the original language. There is no need for special content, but the length of words should match the language.

And after the second paragraph follows the third paragraph. Hello, here is some text without a meaning. This text should show what a printed text will look like at this place. If you read this text, you will get no information. Really? Is there no information? Is there a difference between this text and some nonsense like “Huardest gefburn”? Kjift – not at all! A blind text like this gives you information about the selected font, how the letters are written and an impression of the look. This text should contain all letters of the alphabet and it should be written in of the original language. There is no need for special content, but the length of words should match the language.

Acknowledgments

After this fourth paragraph, we start a new paragraph sequence. Hello, here is some text without a meaning. This text should show what a printed text will look like at this place. If you read this text, you will get no information. Really? Is there no information? Is there a difference between this text and some nonsense like “Huardest gefburn”? Kjift – not at all! A blind text like this gives you information about the selected font, how the letters are written and an impression of the look. This text should contain all letters of the alphabet and it should be written in of the original language. There is no need for special content, but the length of words should match the language.

Hello, here is some text without a meaning. This text should show what a printed text will look like at this place. If you read this text, you will get no information. Really? Is there no information? Is there a difference between this text and some nonsense like “Huardest gefburn”? Kjift – not at all! A blind text like this gives you information about the selected font, how the letters are written and an impression of the look. This text should contain all letters of the alphabet and it should be written in of the original language. There is no need for special content, but the length of words should match the language.

This is the second paragraph. Hello, here is some text without a meaning. This text should show what a printed text will look like at this place. If you read this text, you will get no information. Really? Is there no information? Is there a difference between this text and some nonsense like “Huardest gefburn”? Kjift – not at all! A blind text like this gives you information about the selected font, how the letters are written and an impression of the look. This text should contain all letters of the alphabet and it should be written in of the original language. There is no need for special content, but the length of words should match the language.

Table of Contents

| | |
|--|-----|
| LIST OF TABLES | v |
| LIST OF FIGURES | vi |
| ABSTRACT | vii |
| 1. PROBLEM FORMULATION | 1 |
| 1.1 Video Frame Deletion Detection | 2 |
| 2. PROPOSED APPROACH | 4 |
| 2.1 Prediction Error Sequence Extraction | 4 |
| 2.1.1 Proposed Prediction Error Sequence Extractor for H.264 | 7 |
| 2.2 Proposed Detection Algorithm | 8 |
| A. SOME APPENDIX HEADING | 10 |
| B. ANOTHER APPENDIX HEADING | 12 |
| VITA | 14 |

List of Tables

List of Figures

| | | |
|-----|---|---|
| 1.1 | Generalized System for Frame Deletion Detection | 1 |
| 2.1 | Comparison Between MPEG-2 Detection Methods used on MPEG-2 and H.264 | 6 |
| 2.2 | Comparison Between Prediction Error Sequences and Fingerprint Signals from two Similar Videos | 9 |

Abstract

New Methods for Detecting Frame Deletion in Modern Video

Hunter Kippen

Dr. Matthew Stamm, Ph.D.

And after the second paragraph follows the third paragraph. Hello, here is some text without a meaning. This text should show what a printed text will look like at this place. If you read this text, you will get no information. Really? Is there no information? Is there a difference between this text and some nonsense like “Huardest gefburn”? Kjift – not at all! A blind text like this gives you information about the selected font, how the letters are written and an impression of the look. This text should contain all letters of the alphabet and it should be written in of the original language. There is no need for special content, but the length of words should match the language.

After this fourth paragraph, we start a new paragraph sequence. Hello, here is some text without a meaning. This text should show what a printed text will look like at this place. If you read this text, you will get no information. Really? Is there no information? Is there a difference between this text and some nonsense like “Huardest gefburn”? Kjift – not at all! A blind text like this gives you information about the selected font, how the letters are written and an impression of the look. This text should contain all letters of the alphabet and it should be written in of the original language. There is no need for special content, but the length of words should match the language.

Hello, here is some text without a meaning. This text should show what a printed text will look like at this place. If you read this text, you will get no information. Really? Is there no information? Is there a difference between this text and some nonsense like “Huardest gefburn”? Kjift – not at all! A blind text like this gives you information about the selected font, how the letters are written and an impression of the look. This text should contain all letters of the alphabet and it should be written in of the original language. There is no need for special content, but the length of words should match the language.

Chapter 1: Problem Formulation

Detecting frame deletion in a video requires detecting the structural changes in a video due to the deletion process. In particular, Wang and Farid's work on temporal traces for detecting frame deletion shows that for MPEG-2 video, the P-frame prediction error can be formulated into a sequence, which can be monitored to detect frame deletion. Both Wang and Farid, and Stamm use a system like in Fig. 1.1 to detect frame deletion. The prediction error sequence $e(n)$ is extracted from the decoded video file and processed to produce detection features. Wang and Farid's work did not propose features for automatic detection, and instead relied on visual inspection of the DFT of the prediction error sequence.

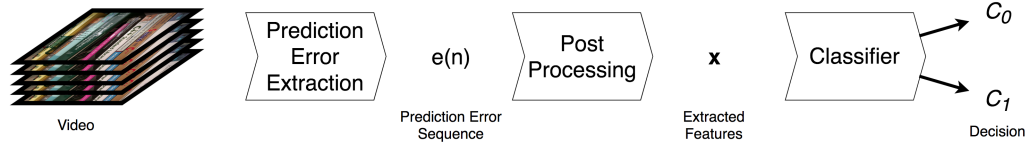


Figure 1.1: Generalized System for Frame Deletion Detection

This generalized system will also work with H.264 encoded video as well. While video encoding has advanced significantly, the fundamental structures of a compressed digital video have remained unchanged. Videos are still stored in GOPs with different numbers of I, P, and B-frames. However, the prediction error extraction and post processing steps must be altered or augmented to remain robust to said advances in video compression.

As this work is concerned particularly with the detection of frame deletion in H.264, we have made the following assumptions. First, we assume that all altered video has undergone re-compression. In fact, since most consumer video recording devices do not have the storage capability or processing power to record high-definition raw video, it is assumed that all video sources have been compressed by either MPEG-4 or H.264, and that all frame deleted video will be re-compressed using H.264,

where the reencoding is set to match the GOP structure of the source video.

In addition, it is assumed that all videos that are passed to the detector are of sufficient length to make a classification. Without multiple full GOPs, the presence of a deletion fingerprint is negligible. Lastly, we make the assumption that if indeed frames have been removed from a video, they have not been removed from the end of the video. The detection features are dependent on differences between the structure of the prediction error sequences in natural videos versus videos with frame deletion. When frames are removed from the end of the video sequence, this difference is not observable.

A user of our proposed system will not need physical access to a specific device to analyze a video captured by the device. The system should accept videos of an arbitrary length, and will not require metadata unrelated to video playback to be in tact. It will work with videos of any resolution, frame rate, or GOP structure. Also, as our approach will be data driven, it is imperative that a user have access to a sufficient database of videos with known labels.

1.1 Video Frame Deletion Detection

Detecting frame deletion is a binary classification problem. Given a Video V , there are two possible classes:

$$\begin{aligned} C_0 &: \text{The video is genuine, and has not had frames removed from it.} \\ C_1 &: \text{The video is altered, and has had frames removed from it.} \end{aligned} \tag{1.1}$$

Note that in this case, *genuine* refers to the fact that the video has not undergone any frame deletion. A video may have underwent other post processing operations such as color correction and re-sizing but not have had any frames removed. In this case, the video would be said to be genuine. From this point forward, any mention of a genuine video simply refers to a video that has not had frames removed from it.

In general, it is difficult to classify based on the entirety of a video directly, thus the problem must be reworked. As shown above, a feature extraction system will be used to produce the P-frame prediction error sequence $e(n)$, and a feature vector \mathbf{x} . The feature vector ideally contains

information about the prediction error sequence that can perfectly separate the two classes. As such, the classification problem is as follows. Given a feature vector \mathbf{x}' , it belongs to one of two classes:

$$\begin{aligned} C_0 : \mathbf{x}' & \text{resulted from a genuine video that has undergone no additional processing.} \\ C_1 : \mathbf{x}' & \text{resulted from an altered video which has had frames removed from it.} \end{aligned} \tag{1.2}$$

In the following chapter, we will propose both a new method for extracting $e(n)$, and additional augmentations to \mathbf{x} that allow for improved separation of data and increased robustness of the overall system.

Chapter 2: Proposed Approach

2.1 Prediction Error Sequence Extraction

In previous work on frame deletion detection in MPEG-2, the prediction error sequence was extracted directly from the video decoder using the DCT coefficients of the prediction error residuals located in the compressed video file. The prediction error was averaged over all macroblocks in a frame. This prediction error was then stored as a sequence. Due to the nature of the correlation between P-frame prediction errors across a single GOP, any prediction made across GOP boundaries would result in increased prediction error. Wang and Farid showed that for fixed GOP video, the increase in average prediction error is periodic with respect to the number of frames deleted from the video. Stamm's work expands the idea of the prediction error trace by introducing the formulation of the problem as detecting the presence of a fingerprint signal $s(n)$. As H.264 uses variable GOP structures we will only be concerned with the model defined for variable GOP video. Stamm defines the model of $s(n)$ as

$$s(n) = \beta \mathbf{1}(\Theta(n) = 0). \quad (2.1)$$

where $\beta > 0$ is a constant and $\Theta(n)$ is a random variable distributed over the set $\{0, 1\}$. This model corresponds to modeling the fingerprint signal as randomly occurring sequence of discrete impulses with a magnitude of β . From this model he poses the detection of frame deletion as distinguishing between two hypotheses:

$$\begin{aligned} H_0 : e(n) &= e_1(n). \\ H_1 : e(n) &= e_2(n) = e_1(n) + s(n)e_1(n). \end{aligned} \quad (2.2)$$

Thus, detection of frame deletion is detection of the presence of the modulated fingerprint signal $s(n)e_1(n)$. Given an unknown video, Stamm first makes an approximation of the unaltered P-frame

prediction error sequence. To do this he uses a median filter with a filter width of 3.

$$\hat{e}(n) = \text{median}\{e(n-1), e(n), e(n+1)\}. \quad (2.3)$$

Thus, the relationship between the estimate and $e_1(n)$ is

$$e_1(n) = \hat{e}(n) + \epsilon(n). \quad (2.4)$$

where $\epsilon(n)$ is a zero mean random variable representing estimation error.

Using the estimate of the unaltered P-frame prediction error sequence, Stamm calculates $\hat{s}(n)$, which is an estimate of the fingerprint signal modulated by the prediction error sequence as defined by

$$\hat{s}(n) = \max(e(n) - \hat{e}(n), 0). \quad (2.5)$$

The estimate of the fingerprint signal is floored at 0, as the model of the $s(n)$ dictates that it must be greater than or equal to 0. This estimate of the fingerprint signal can be used to build a detector. The decision function found by Stamm for variable GOP video is

$$\delta_{var} = \begin{cases} H_0 & \text{if } \frac{1}{N} \sum_{n=1}^N |\hat{s}(n)| < \tau_{var} \\ H_1 & \text{if } \frac{1}{N} \sum_{n=1}^N |\hat{s}(n)| \geq \tau_{var} \end{cases} \quad (2.6)$$

In MPEG-2, a P-frame is encoded by searching the previous anchor frame for the macroblock which incurs the least error. This means that the average prediction error for a single P-frame is only associated with the previous I or P-frame. H.264 expands the capabilities of its motion compensation and estimation system by allowing prediction from multiple previous frames (and subsequent frames in the case of B-frames). If the prediction error trace is extracted via the codec for H.264, the average prediction error associated with one frame is comprised of a linear combination of the average prediction error associated with motion vectors that map to the different anchor frames

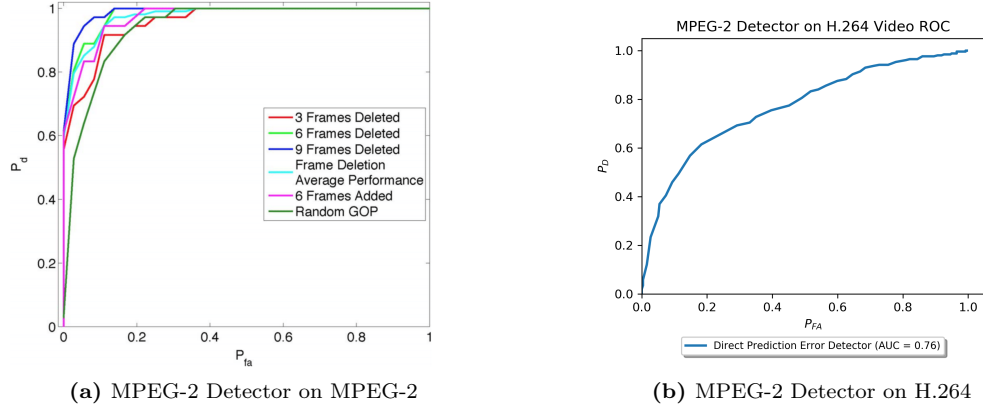


Figure 2.1: Comparison Between MPEG-2 Detection Methods used on MPEG-2 and H.264

used in the motion estimation and compensation process. In H.264, the prediction error sequence $e(n)$ can be defined. Thus, cross GOP predictions are smoothed out in such a way that it makes the fingerprint energy detector in Stamm's paper perform inadequately.

To test this, we collected 230 videos from a cell phone camera (the ASUS ZenFone 3 Laser), and generated an altered video with 15 frames removed from the beginning corresponding to each collected video. The encoding parameters were kept constant, and we fixed the GOP of the altered videos to that of the unaltered videos. In this particular case, the GOP structure was 30 frames in length, with 1 I-frame followed by 29 P-frames. We extracted the prediction error sequence directly from the codec, and measured the estimated fingerprint energy as described above.

As shown in Fig. 2.1 the probability of detection using the methodology from MPEG-2 videos does not translate to H.264, particularly at low false alarm rates. This is due to a limitation in the model used by Stamm above in equation 2.1. As it is possible to predict across multiple previous anchor frames in H.264, the contribution of the fingerprint signal is variable over time. This variation is also not regular, as scene content and motion determine how many cross GOP predictions are present in each P-frame. Thus we propose the updated model for the fingerprint signal as

$$s(n) = \beta(n) \mathbf{1}(\Theta(n) = 0). \quad (2.7)$$

where $\beta(n)$ is now a random variable that takes values in $\mathbb{R}_{\geq 0}$. Thus, we propose the following methodology for extracting the prediction error sequence in H.264.

2.1.1 Proposed Prediction Error Sequence Extractor for H.264

The goal of the proposed extraction algorithm is to maximize the probability that should frame deletion exist, a given measurement of the prediction error comes from a cross-GOP prediction. To this end, instead of directly measuring the prediction error from the DCT coefficients from the decoder, we decode the frame of interest and store the motion vectors associated with said frame. For each macroblock in the current frame, we use its associated motion vector to find the x and y coordinates defining the source macroblock which provides the least error mapping from a previous anchor frame. Then for a previous anchor frame, we subtract the pixels in the source macroblock from the destination macroblock in the current frame. This leaves us with a prediction error residual associated with the macroblock. We then calculate the average absolute value of this residual.

We repeat this process for each of D previous anchor frames. Then we store these prediction error values in a matrix M , where each row of the matrix corresponds to the errors associated with a single motion vector, and the columns are the errors associated with each of the previous anchor frames.

After obtaining the M matrix for the current frame, we create the matrix \tilde{M} defined like so:

$$\tilde{M}_{i,j} = \mathbf{1} \left(j = \underset{l}{\operatorname{argmin}} (M_{i,l}) \right) * M_{i,j} \quad (2.8)$$

Thus \tilde{M} is a copy of M where the non-zero entries in each row correspond to the minimum average error associated with the macroblock and all other elements are zero. Effectively, the column index of the non-zero entry is an estimation of which previous frame the motion vector associated with the macroblock maps to in the decoding process.

Further processing is done on \tilde{M} to output only a single prediction error value. First, \tilde{M} is

reduced into a vector P , such that:

$$P_j = \frac{1}{N_j} \sum_i \tilde{M}_{i,j} \quad (2.9)$$

Where N_j is the number of non-zero elements in the j^{th} column of \tilde{M} . Then, the reported prediction error e is set as the maximum value of P . This entire process is repeated for every P frame in the video. This method of error extraction estimates which previous anchor frame contributes the maximum error per macroblock to the overall prediction error residual obtained by the codec for a given P frame. Since prediction across GOP boundaries results in spikes in the prediction error, the anchor frame that contributes the most error is most likely to be from a different original GOP. In this manner, we obtain a trace that is resilient to advances made in the motion compensation and estimation process in modern codecs as well as robust to variable frame rates and dynamic GOP structures.

2.2 Proposed Detection Algorithm

In addition to the new methods for prediction error extraction, we propose an expanded detection algorithm to better capture the statistical differences between videos. In fact, depending on scene content, video capture settings, and the amount of motion captured in a single recording, the prediction error sequence and fingerprint signal exhibit different structural behavior. This is true even for videos captured from a single camera model. Figure 2.2 shows this clearly. The two videos were captured from an LG Nexus 5X using the high quality 1080p capture mode. Both videos were shot using similar scene content, but the amount of motion in each video is different. The first video was shot with high motion, while the second video was comparatively low motion. The top row shows the different prediction error sequences, while the bottom row shows the different estimated fingerprint signals.

Thus there is a need to construct a parametric model of the fingerprint signal and the prediction error sequence to capture these statistical variations. We propose constructing an autoregressive (AR) model of both the fingerprint signal and the prediction error sequence. Over the time period

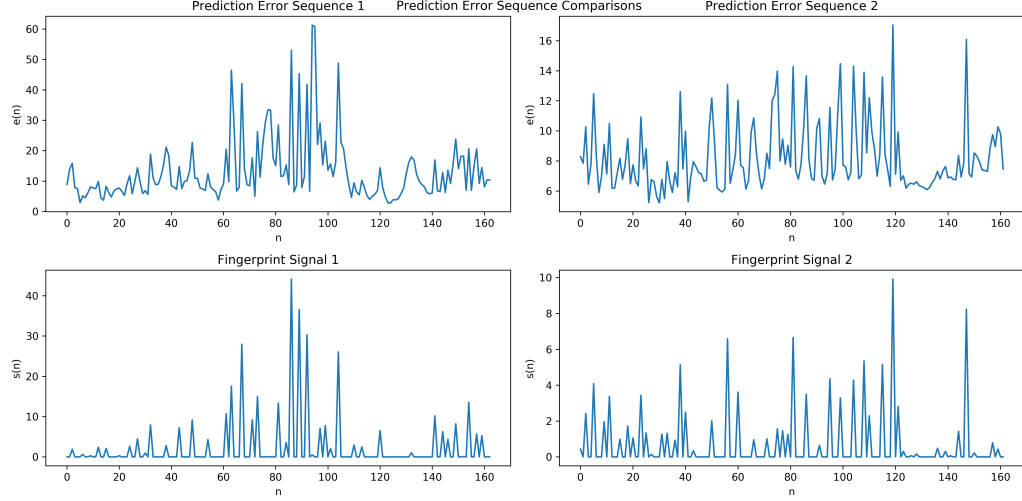


Figure 2.2: Comparison Between Prediction Error Sequences and Fingerprint Signals from two Similar Videos

of one GOP it can be said that the signals are mostly stationary. The parameters of the AR models are added to a feature vector along with the fingerprint energy described above. In order to capture the degree to which this is true, the error variance of each AR model is also included. In addition, some basic statistical features are included to scale the overall decision surface. We propose including the mean and variance of both the prediction error sequence and fingerprint signal to the feature vector as well.

Note that after creating a feature vector for each video, the feature vector is quite large. It is inadvisable to create a probabilistic model of the feature vector for classification. Instead, we propose using a discriminative function to map an incoming feature vector directly to the set of natural videos or the set of videos altered by frame deletion. As such, we propose using a Support Vector Machine (SVM) classifier with a Radial Basis Kernel function for classification.

Appendix A: Some Appendix Heading

This is the second paragraph. Hello, here is some text without a meaning. This text should show what a printed text will look like at this place. If you read this text, you will get no information. Really? Is there no information? Is there a difference between this text and some nonsense like “Huardest gefburn”? Kjift – not at all! A blind text like this gives you information about the selected font, how the letters are written and an impression of the look. This text should contain all letters of the alphabet and it should be written in of the original language. There is no need for special content, but the length of words should match the language.

And after the second paragraph follows the third paragraph. Hello, here is some text without a meaning. This text should show what a printed text will look like at this place. If you read this text, you will get no information. Really? Is there no information? Is there a difference between this text and some nonsense like “Huardest gefburn”? Kjift – not at all! A blind text like this gives you information about the selected font, how the letters are written and an impression of the look. This text should contain all letters of the alphabet and it should be written in of the original language. There is no need for special content, but the length of words should match the language.

After this fourth paragraph, we start a new paragraph sequence. Hello, here is some text without a meaning. This text should show what a printed text will look like at this place. If you read this text, you will get no information. Really? Is there no information? Is there a difference between this text and some nonsense like “Huardest gefburn”? Kjift – not at all! A blind text like this gives you information about the selected font, how the letters are written and an impression of the look. This text should contain all letters of the alphabet and it should be written in of the original language. There is no need for special content, but the length of words should match the language.

Hello, here is some text without a meaning. This text should show what a printed text will look like at this place. If you read this text, you will get no information. Really? Is there no information? Is there a difference between this text and some nonsense like “Huardest gefburn”? Kjift – not at

all! A blind text like this gives you information about the selected font, how the letters are written and an impression of the look. This text should contain all letters of the alphabet and it should be written in of the original language. There is no need for special content, but the length of words should match the language.

This is the second paragraph. Hello, here is some text without a meaning. This text should show what a printed text will look like at this place. If you read this text, you will get no information. Really? Is there no information? Is there a difference between this text and some nonsense like “Huardest gefburn”? Kjift – not at all! A blind text like this gives you information about the selected font, how the letters are written and an impression of the look. This text should contain all letters of the alphabet and it should be written in of the original language. There is no need for special content, but the length of words should match the language.

Appendix B: Another Appendix Heading

And after the second paragraph follows the third paragraph. Hello, here is some text without a meaning. This text should show what a printed text will look like at this place. If you read this text, you will get no information. Really? Is there no information? Is there a difference between this text and some nonsense like “Huardest gefburn”? Kjift – not at all! A blind text like this gives you information about the selected font, how the letters are written and an impression of the look. This text should contain all letters of the alphabet and it should be written in of the original language. There is no need for special content, but the length of words should match the language.

After this fourth paragraph, we start a new paragraph sequence. Hello, here is some text without a meaning. This text should show what a printed text will look like at this place. If you read this text, you will get no information. Really? Is there no information? Is there a difference between this text and some nonsense like “Huardest gefburn”? Kjift – not at all! A blind text like this gives you information about the selected font, how the letters are written and an impression of the look. This text should contain all letters of the alphabet and it should be written in of the original language. There is no need for special content, but the length of words should match the language.

Hello, here is some text without a meaning. This text should show what a printed text will look like at this place. If you read this text, you will get no information. Really? Is there no information? Is there a difference between this text and some nonsense like “Huardest gefburn”? Kjift – not at all! A blind text like this gives you information about the selected font, how the letters are written and an impression of the look. This text should contain all letters of the alphabet and it should be written in of the original language. There is no need for special content, but the length of words should match the language.

This is the second paragraph. Hello, here is some text without a meaning. This text should show what a printed text will look like at this place. If you read this text, you will get no information. Really? Is there no information? Is there a difference between this text and some nonsense like

“Huardest gefburn”? Kjift – not at all! A blind text like this gives you information about the selected font, how the letters are written and an impression of the look. This text should contain all letters of the alphabet and it should be written in of the original language. There is no need for special content, but the length of words should match the language.

And after the second paragraph follows the third paragraph. Hello, here is some text without a meaning. This text should show what a printed text will look like at this place. If you read this text, you will get no information. Really? Is there no information? Is there a difference between this text and some nonsense like “Huardest gefburn”? Kjift – not at all! A blind text like this gives you information about the selected font, how the letters are written and an impression of the look. This text should contain all letters of the alphabet and it should be written in of the original language. There is no need for special content, but the length of words should match the language.

Vita

After this fourth paragraph, we start a new paragraph sequence. Hello, here is some text without a meaning. This text should show what a printed text will look like at this place. If you read this text, you will get no information. Really? Is there no information? Is there a difference between this text and some nonsense like “Huardest gefburn”? Kjift – not at all! A blind text like this gives you information about the selected font, how the letters are written and an impression of the look. This text should contain all letters of the alphabet and it should be written in of the original language. There is no need for special content, but the length of words should match the language.

Hello, here is some text without a meaning. This text should show what a printed text will look like at this place. If you read this text, you will get no information. Really? Is there no information? Is there a difference between this text and some nonsense like “Huardest gefburn”? Kjift – not at all! A blind text like this gives you information about the selected font, how the letters are written and an impression of the look. This text should contain all letters of the alphabet and it should be written in of the original language. There is no need for special content, but the length of words should match the language.

This is the second paragraph. Hello, here is some text without a meaning. This text should show what a printed text will look like at this place. If you read this text, you will get no information. Really? Is there no information? Is there a difference between this text and some nonsense like “Huardest gefburn”? Kjift – not at all! A blind text like this gives you information about the selected font, how the letters are written and an impression of the look. This text should contain all letters of the alphabet and it should be written in of the original language. There is no need for special content, but the length of words should match the language.

