

**FAO Dr. Amina Souag, Dr. Hannan  
Azhar**

**???????**

**BSc (Hons.) Computer Science**

**2023-2024**

**Individual Project 40**

**Title: ??????**

**Author: Jamie Pinnington**

**Supervisors: Dr. Amina Souag, Dr. Hannan Azhar**

**Email: JP878@canterbury.ac.uk**

This report is submitted in partial fulfilment of the requirement for the  
BSc in *Computer Science* at Canterbury Christ Church University

I declare that this report is my own original work containing no personal data as defined in the Data Protection Act (1998) and that I have read, understood and accept the University's regulations on plagiarism/intellectual property rights/research ethics (in particular the Research Governance Handbook) and the IP 40 Module Handbook.

Further, I accept that digital and/or hard copies of my Individual Project 40, or parts thereof, may be made available to other students, individuals and organisations after it has been marked.

Finally, I accept that no copy of my Individual Project 40 will ever be returned regardless of the circumstances.

Signed *Jamie Pinnington*

Date of Submission: ??????

## 1. Acknowledgements

## Contents

<b>1 Acknowledgements</b>	<b>1</b>
<b>2 Introduction:</b>	<b>6</b>
<b>3 Literature Review:</b>	<b>6</b>
3.1 Introduction . . . . .	6
3.1.1 Background . . . . .	6
3.1.2 Purpose . . . . .	6
3.1.3 Scope . . . . .	7
3.1.4 Research Questions . . . . .	7
3.2 Methodology . . . . .	8
3.2.1 Literature Identification . . . . .	8
3.2.2 Literature Evaluation . . . . .	8
3.3 Historical Context (RQ-7) . . . . .	9
3.4 Methods, Tools, and Techniques (RQ-1, RQ-2, RQ-3) . . . . .	15
3.4.1 Image-based Methods . . . . .	15
3.4.2 Video-based Methods . . . . .	15
3.4.3 Framework-based Methods and Tools . . . . .	15
3.5 Comparative Analysis of Machine Learning Models (RQ-1) . . . . .	16
3.6 Challenges in ASL Fingerspelling Recognition (RQ-5, RQ-6) . . . . .	17
3.7 Overcoming Obstacles (RQ-5) . . . . .	17
3.8 State of the Art and Real-World Applications (RQ-4, RQ-7) . . . . .	17
3.9 Future Directions and Open Challenges (RQ-7) . . . . .	18
3.10 Ethical and Societal Considerations . . . . .	18
3.11 Conclusion . . . . .	18
<b>4 MAIN CHAPTERS:</b>	<b>22</b>
4.1 Chapter 1: . . . . .	22
<b>5 DEVELOPMENT NOTES</b>	<b>22</b>
5.1 Algorithm Pseudocode . . . . .	22
5.2 Model Art . . . . .	22
<b>Appendix A Glossary</b>	<b>A1</b>

<b>Appendix B</b>	<b>Marking Scheme</b>	<b>B1</b>
<b>Appendix C</b>	<b>Changes to the Project Initiation Document</b>	<b>C1</b>
<b>Appendix D</b>	<b>Current Environment Investigation Report</b>	<b>D1</b>
<b>Appendix E</b>	<b>Requirements Specification</b>	<b>E1</b>
<b>Appendix F</b>	<b>Design Report</b>	<b>F1</b>
<b>Appendix G</b>	<b>Implementation</b>	<b>G1</b>
<b>Appendix H</b>	<b>Testing</b>	<b>H1</b>
<b>Appendix I</b>	<b>User Guide</b>	<b>I1</b>
<b>Appendix J</b>	<b>Project Management</b>	<b>J1</b>
<b>Appendix K</b>	<b>Meetings With Supervisor</b>	<b>K1</b>
<b>Appendix L</b>	<b>Agile Development: Timebox 1</b>	<b>L1</b>
<b>Appendix M</b>	<b>Agile Development: Timebox 2</b>	<b>M1</b>
<b>Appendix N</b>	<b>Agile Development: Timebox 3</b>	<b>N1</b>

**List of Figures**

J.1	Work Breakdown Structure: Develop Model . . . . .	J1
J.2	Work Breakdown Structure: Build Application . . . . .	J2
J.3	Work Breakdown Structure . . . . .	J3
J.4	Gantt Chart . . . . .	J4
J.5	Risk register . . . . .	J5

**List of Tables**

1	Summary and Analysis of ASL Fingerspelling Recognition Models (2018-2023)	12
---	---	----

## 2. Introduction:

**T**<sup>HIS</sup>

## 3. Literature Review:

### 3.1. Introduction

#### 3.1.1. Background

- why is this topic (ASL) important? (e.g., accessibility, communication)
- what is the specific problem being addressed? (e.g., fingerspelling recognition) (bridges the gap in communication and enhances the learning/usage of ASL.) (makes ai more accessible to this audience?) ASL ,
- what is the impact of an AI recognizer for this problem? (e.g., enables real-time communication, improves accessibility, etc.) (can this effect broader technology?)

Sign language is the primary form of communication for the deaf and hard of hearing community. It allows communication when the spoken language is not possible, and or when the speaker or receiver is deaf or hard of hearing. Depending on the situation, and like any language, it requires both parties to be fluent in the language to communicate effectively. However, this is not always the case. American Sign Language(ASL) is a complete, complex language that employs signs made with the hands and other movements, including facial expressions and postures of the body, and is used natively in the United States of America and globally by many individuals.

Whilst no attempt has officially been made to survey the language, and most current estimates are based off of historical surveys that prove to be inaccurate Mitchell et al. (2006). It is estimated that there are over 1 million signers Ethnologue (2023), but others estimates are as high as 2 million Mitchell et al. (2006). ASL communicates through a variety of means including gestures, non-manual markers and lexical signs. The most understood are lexical vocabulary, each corresponding to a word or morpheme. Gestures and non-manual markers such as facial expression can complement and convey more interactive or meaningful lexical signs. Additional constructs include usage of space, role shifting and classifiers.

#### 3.1.2. Purpose

- what is the purpose of the review?
- (e.g., to identify the state of the art in ASL fingerspelling recognition)

- (to identify the challenges and opportunities in ASL fingerspelling recognition)
- (to identify the most promising techniques for ASL fingerspelling recognition)
- \* primary purpose is to build our own model, but we need to know what's out there first. \*

### 3.1.3. *Scope*

- what is the scope of the review? (e.g., ASL fingerspelling recognition) (what is the scope of the problem? (e.g., real-time recognition of fingerspelling gestures) (what is the scope of the solution? (e.g., image-based recognition of fingerspelling gestures) (what is the scope of the evaluation? (e.g., accuracy, speed, etc.)
- what is the scope of the literature? (e.g., papers published in the last 5 years) (what is the scope of the sources? (e.g., peer-reviewed journal articles, conference papers, etc.)
- what we're not covering.
- only recognition and translation of ASL \*fingerspelling\* (not full ASL).
- specific the application/methodology (e.g., video-based recognition of fingerspelling gestures) (live/stream???)

### 3.1.4. *Research Questions*

- RQ-1: Comparative Analysis of Machine Learning Models: What are the strengths and weaknesses of different machine learning (ML) models, such as Convolutional Neural Networks (CNNs), Recurrent Neural Networks (RNNs), and Transformer models, in the context of ASL fingerspelling recognition?
- RQ-2: Performance Evaluation: How do various machine learning models perform in terms of accuracy, processing speed, and reliability for ASL fingerspelling recognition under different conditions (e.g., varying lighting, hand positions, backgrounds)?
- RQ-3: Dataset and Model Suitability: How does the choice of dataset, including its size, diversity, and quality, influence the effectiveness of different machine learning models in recognizing ASL fingerspelling?
- RQ-4: Real-World Applications: Considering practical applications like kiosk systems, which machine learning models offer the best balance between technical performance and user experience for ASL fingerspelling recognition?



- RQ-5: Technical Challenges: What technical challenges are commonly faced across different machine learning models in ASL fingerspelling recognition, and how adaptable are these models to address such challenges?
- RQ-6: Impact of Environment Variables: To what extent do environmental variables (like hand orientation, motion speed, and background noise) affect the performance of different machine learning models in ASL fingerspelling recognition?
- RQ-7: State of the Art and future directions: What are the most recent and influential works in the field of ASL fingerspelling recognition, and what are the emerging trends and future directions?

### 3.2. Methodology

#### 3.2.1. Literature Identification

This literature search was completed using the databases IEEE Xplore, Google Scholar, ACM Digital Library, and ScienceDirect. Search terms such as "ASL fingerspelling recognition", "Deep learning for ASL recognition", "ASL recognition with CNN", "Accuracy of ASL recognition models", "Latest trends in ASL recognition", and "ASL recognition in real-time" were used alone and in conjunction with boolean operators "AND", "OR" to refine the search results. The search was limited to papers published in the last 5 years, and only peer-reviewed journal articles, conference papers, and high quality theses were considered. The search was also limited to papers written in English. The search was conducted in between November and December 2023, and the results were filtered to include only papers that were published between 2018 and 2023. In addition to this search which was completed in order to find relevant literature for the model architecture comparison, the references of the papers that were selected were used to find additional relevant literature, that dates further back in order to substantiate our historical context and technical background.

#### 3.2.2. Literature Evaluation

Of the literature that fit out search criteria, the most relevant papers were selected based on the following criteria: the relevance to the research questions, papers that specifically address ASL fingerspelling, ML models in sign language interpretation, papers that used widely recognized datasets relevant to ASL recognition, editorials, opinion pieces, and non-peer reviewed articles were excluded. Papers that were preferred had a clear methodology, defined objectives and analysis of data. Papers with high citation counts were also given preference. Papers that were selected were read in full, and the results were summarized in a table, which is included in the results section.

- By understanding, the insight gained from this review will be used to inform the design and implementation of our own model.
- \* evidence based approach is necessary to be impactful \*

### 3.3. *Historical Context (RQ-7)*

- Provide a brief overview of the evolution of ASL fingerspelling recognition.
- Highlight key milestones and breakthroughs in the field's history.
- Connect historical developments to current trends and future directions.

Early approaches to sign language recognition, according to Saeed et al. (2022) used robotic like data/power gloves which were wired, with sensors to capture hand movements and gestures. They aimed to record the finger position and flexion in order to classify shapes. These approaches were limited by the need for specialized hardware and the inability to capture facial expressions and other non-manual markers. Rule-based classifiers did the legwork by detecting specific input pattern of sensors to an output by programmatic rules. This approach was not practical or user-friendly.

The move to ML was occurring in parallel to hardware based approaches, as vision-based approaches were developed to overcome the limitations of the hardware based approaches and was instrumental in the development of sign language recognition. This approach used computer vision techniques to detect and track the hand and fingers. They were able to capture more information than the hardware based approaches, but were limited by the need for a controlled environment and the inability to capture facial expressions and other non-manual markers, just as the hardware based approaches were. At this stage, there were no large datasets developed, and the datasets that were available were not standardized, and were not publicly available. The vocabulary was relatively small von Agris et al. (2008), which meant that the recognition was limited to a small number of signs. The recognition was also limited to a single signer, and was not robust to variations in lighting, hand orientation, and background noise.

A great deal of focus was on feature extraction and classification, various algorithms were being pursued to extract hand features like posture. Hidden Markov Models (HMMs) were used to solve this temporal sequential task, where each sign or gesture is defined by the transition from one state to another. HMMs use the transition state to understand meaning. This approach was limited by the need for a large amount of data to train the model von Agris et al. (2008).

Another leap with vision was the usage of support vector machines (SVMs) together with HMMs to enhance classification. SVMs were more effective at classifying spatial features such as

hand shape and geolocation of digits, and videos that have depth, where gestures or shapes could look similar in 2D or 3D Vogler and Metaxas (1999).

Moving around the late 00's, a significant feat was neural networks such as used in Munib et al. (2007), which used a 3-layer network with backpropagation and Hough transform. Although 92.3% accuracy was achieved, this was still comparable to models using SVMs and HMMs, as well as the hardware based approaches before. The dataset was still limiting, with "300 samples of hand sign images; 15 images for each sign." Munib et al. (2007).

Perhaps the greatest leap was the rise of deep learning, as neural networks got deeper in terms of model layering. Image and video processing research as a whole was in full swing, convolutional neural networks (CNNs), recurrent neural networks (RNNs), and Long short-term memory (LSTMs), were a few which had significant impact. CNNs were adept at automatically extracting and learning

Yann Lecun is credited with setting the precedence of CNNs in 1998, with the LeNet-5 architecture Lecun et al. (1998). This was a significant leap in the field of computer vision, and was the first time that a model was able to learn features automatically, which is why due to greatly increased GPU processing power, CNNs came back into the fold at the 2012 ImageNet challenge Krizhevsky et al. (2012). CNNs are ideal and particularly adept at processing data that is grid-like, as in images that have dimensions. A series of layers are used to extract and identify features by breaking down the image into smaller parts, understanding that, and over and over, and combining them to understand the whole image. Functions and pooling layers are used to optimize the output for classification.

RNNs are a type of neural network that are adept at processing sequential data, such as text, audio, and video. They are able to remember previous inputs and use that information recurrently, because the network has a directed cycle network, meaning information can persist inside the network Sherstinsky (2020). This is particularly useful for ASL recognition, as the signs are sequential, and the order of the signs is important.

LSTMs are a type of RNN that are able to remember information for long periods of time, and are able to overcome the vanishing gradient problem that is common in RNNs Sherstinsky (2020).

While these models and approaches are valid for the challenge of American Sign Language, commonly fingerspelling isn't factored in, and these models struggle to perform well on fingerspelling. This may be because fingerspelling is a sequential task where the order of the letters is important, and that the subtle differences between letters are difficult to distinguish. More so, in practise the signed spellings change rapidly and do not necessarily finish a movement, the hand is continuously transitioning sequentially from one letter to the next to form a word. This is a challenge for models that are not able to capture the temporal nature of the task.

Currently, the experimental methods are Connectionist Temporal Classification (CTC) Graves

et al. (2006); Shi et al. (2018), Attention Bahdanau et al. (2016), Transformers Vaswani et al. (2023), and using large language models (LLMs) to improve accuracy among others.

Table 1: Summary and Analysis of ASL Fingerspelling Recognition Models (2018-2023)

Reference	Model Used	Framework	Dataset	Key Findings	Performance Metrics	Challenges Addressed
S Kumar et al. (2018)	RNN, LSTM, Attention, Encoder/Decoder	[Not Specified]	NCSLGR Corpus	Recognition and translation of ASL glosses	GRR: 86%, GER: 23%	Real-time recognition and translation
Weerasooriya and Ambegoda (2022)	RF, KNN, LR	[Not specified]	FASSL custom dataset	Developed a classifier for static signs using a small dataset	Accuracy: 87.9% (correct estimates)	Pose classification with limited data
Cihan Camgoz et al. (2020)	Transformers with CTC loss	PyTorch	PHOENIX14T	State-of-the-art results in recognition and translation	WER, BLEU-4 scores	Translation from sign language videos to spoken language sentences
Abiyev et al. (2020)	CNN, SSD, FCN	[Not specified]	Kaggle ASL Fingerspelling	High accuracy, vision-based translation	Accuracy: 92.21%	Real-time translation, robustness in ASL recognition
Bantupalli and Xie (2018)	CNN, LSTM, RNN	OpenCV	Self-created Dataset	Effective recognition with custom CNN model	Accuracy: 98.11%	Robust recognition in controlled environments
Kabade et al. (2023)	ResNet, Bi-LSTM, CTC, Attention	[Not specified]	ChicagoFSWild	Recognition using optical flow and attention, preprocessing for occlusions	Letter accuracy: 57%	Recognition in 'wild' conditions, occlusions

Reference	Model Used	Framework	Dataset	Key Findings	Performance Metrics	Challenges Addressed
Shi et al. (2018)	CNN, LSTM, CTC	Faster R-CNN	Custom YouTube Dataset	Improved accuracy with hand detection	Test Acc: 41.9% with CTC	Recognition in the wild, varying conditions
Shi et al. (2019)	CNN, RNN, CTC, Attention	TensorFlow	ChicagoFSWild, ChicagoFSWild+	Enhanced recognition in uncontrolled environments	Word Error Rate: 27.2	Recognition in diverse and challenging real-world scenarios
Shi et al. (2021)	2D/3D-CNN, Bi-LSTM	OpenPose	ChicagoFSWild, ChicagoFSWild+	Superior detection in uncontrolled environments	AP@IoU: 0.495, MSA: 0.386	Handling fine-grained handshapes and signer's pose
Nguyen and Do (2019)	1) LBP, HOG descriptors, multi-class SVM, 2) End-to-end CNN 3) CNN weights as feature extractor for Linear-kernel SVM	[Not specified]	Massey Dataset	Three diverse methods for fingerspelling recognition	Recognition rate: 97.49%, 98.23%, 98.30%	Adaptability in feature extraction and classification approaches
Chong and Lee (2018)	SVM and DNN	TensorFlow, Scikit-learn	Self-created Dataset	Comparison of SVM and DNN for ASL recognition; effective use of LOO approach for bias avoidance	Recognition rate: 72.79%, 88.79%	Multi-class classification with 36 classes (26 letters and 10 digits)
Bantupalli and Xie (2018)	CNN (Inception) for spatial features, LSTM for temporal features	TensorFlow, Keras	American Sign Language Dataset	Efficient extraction of temporal and spatial features; use of Inception and LSTM models	Accuracy up to 93% (Softmax Layer), 58% (Pool Layer)	Managing longer sequences with LSTM; preventing overfitting with dropout

Reference	Model Used	Framework	Dataset	Key Findings	Performance Metrics	Challenges Addressed
Shi et al. (2022)	FSS-Net (End-to-End Model for Fingerspelling Detection and Text Matching)	[Not Specified]	ChicagoFSWild, ChicagoFSWild+	Introduced explicit temporal localization for fingerspelling search and retrieval. Demonstrated effective fingerspelling detection in varying conditions.	mAP: 0.684 (YouTube), 0.584 (DeafVIDEO), 0.629 (Misc)	Fingerspelling detection in diverse visual conditions; handling open vocabulary and arbitrary-length queries; confusion between similar handshapes; detection failures.
Gajurel et al. (2021)	Fine-Grained Visual Attention with Transformer Model (CTC, CNN, LSTM)	[Not specified]	ChicagoFSWild	Significantly improved state-of-the-art performance in fingerspelling recognition using Transformer-based contextual attention mechanism	Letter Accuracy: 46.96% (dev), 48.36% (test)	Addressed challenges in capturing fine-grained details in unsegmented continuous video data. Focused on improving generalization and regularization of the model.

### 3.4. *Methods, Tools, and Techniques (RQ-1, RQ-2, RQ-3)*

- Offer an overview of the various methods used in ASL fingerspelling recognition.
- Discuss Image/Video-based methods, Framework-based approaches (e.g., MediaPipe), and Hybrid methods in detail.
- Include a discussion of common evaluation metrics and their relevance to different models and conditions.

#### 3.4.1. *Image-based Methods*

The system performs static feature extraction on each individual image. Information from the image such as position of hand, texture, shape, colour, and other features are extracted and used to classify the image. The system doesn't identify or understand the temporal nature of the task, and is limited to the information that can be extracted from the image. ML models such as CNNs and SVMs are used to classify the extracted features to specific letters. This approach wouldn't be suitable for real-time recognition.

#### 3.4.2. *Video-based Methods*

This time, the system performs sequential feature extraction, by extracting temporal features the models can understand the transitions between letters. This is critical for differentiating letters that might be similar when static, or depending on the orientation. For instance, the letter "h" and "u" have the same hand shape, but differ in meaning depending on orientation. Models typically include RNNs and LSTMs which are particularly useful at retaining information over time, which are important for a real-time sequential task, that must understand the order of the letters.

#### 3.4.3. *Framework-based Methods and Tools*

The MediaPipe framework, introduced by Lugaresi et al. (2019), is an open-source collection of libraries and tools designed to facilitate the development and deployment of artificial intelligence (AI) and machine learning (ML) applications. Its implementation is particularly beneficial for creating ML pipelines, offering a suite of pre-trained models that excel in tasks such as gesture recognition and hand segmentation. As a relatively new tool, MediaPipe is gaining traction among developers who require robust solutions for specific layers of the ML pipeline without the need to develop models from scratch.

OpenCV, as detailed by Culjak et al., is a comprehensive open-source library of optimized algorithms that provides extensive support for image and video processing tasks. Widely adopted in the field of computer vision, OpenCV is commonly utilized for detecting signs, numerals, alphabets,



and more, often serving as a cornerstone in the preprocessing stages of machine learning models Srinivasan et al. (2023). Its versatility and performance make it a popular choice for researchers and practitioners working on sign language recognition and related areas.

TensorFlow Abadi et al. (2016) and PyTorch Paszke et al. (2019) are both development systems, each individual ecosystems in their own right where developers can build and train models at scale. They are both open-source, and are both widely used in the field of machine learning. TensorFlow is developed by Google, and PyTorch is developed by Meta. TensorFlow is more mature, and has a larger community, and is more widely used in production. PyTorch is more flexible, and is more popular for research and experimentation. Both frameworks are used in the development of ASL recognition models.

### 3.5. *Comparative Analysis of Machine Learning Models (RQ-1)*

S Kumar et al. (2018) and several others utilize RNN, LSTM, and Attention Mechanisms, highlighting the importance of sequential data processing in sign language. CNN combined with LSTM, as in Shi et al. (2018), indicates the effectiveness of capturing both spatial and temporal features. The use of Transformers, such as in Cihan Camgoz et al. (2020), showcases advanced capabilities in translation tasks.

The study by Weerasooriya and Ambegoda (2022) using RF, KNN, and LR represents traditional machine learning approaches, effective for smaller datasets. In contrast, Chong's comparison Chong and Lee (2018) between SVM and DNN illustrates the evolving landscape from classical to modern neural network-based approaches. Performance Metrics and Dataset Dependency:

High accuracy in controlled environments, like Bantupalli's 98.11% accuracy Bantupalli and Xie (2018), contrasts with moderate success in 'wild' conditions, such as Kabade's 57% letter accuracy Kabade et al. (2023). The choice of datasets, ranging from custom ones to larger, more diverse datasets like PHOENIX14T or ChicagoFSWild, influences the model selection and performance, as seen across multiple studies.

Real-time recognition and translation needs, addressed in studies like Abiyev et al. (2020), demand fast and efficient models. Recognition in uncontrolled environments, as explored by Shi et al. (2019), requires robust models capable of handling diverse and challenging scenarios.

Unique approaches like Shi's FSS-Net (Shi et al., 2022) emphasize innovation in addressing specific challenges like temporal localization and open vocabulary in fingerspelling detection. Gajurel's study (Gajurel et al., 2021) using a Transformer model with fine-grained visual attention highlights efforts in improving model generalization and handling unsegmented continuous video data.

### 3.6. *Challenges in ASL Fingerspelling Recognition (RQ-5, RQ-6)*

- Identify and discuss the technical, data, and real-world challenges in ASL fingerspelling recognition.
- Include a discussion on how environmental variables, such as hand orientation and background noise, affect model performance.

There are currently a numerous amount of challenges facing all types of ML models in ASL fingerspelling recognition. These challenges can be categorized into three main categories: technical, data, and real-world challenges. The variability in hand shape and motion neccitate deeper models in order to learn the complex patterns of large amounts of data Gajurel et al. (2021). Fluent signers can spell quickly and smoothly than novice signers, which can be difficult for models to capture the fluidity of the motion Gajurel et al. (2021). Depending on the angle or position of the hand, occlusion from other fingers or the body can occur, which can make it difficult for models to distinguish between letters for signers like "A" and "S" Shi et al. (2018). Classification is harder because the model cannot see the whole hand, and the model must learn to recognize the letter from a partial view of the hand Shi et al. (2018). Overlapping is another challenge, where the movement from one letter to another too quickly, may overlap as there is no distinct pause, or boundary between the spelt letters. Diverse background and lighting conditions. One of the first steps in any sign language task, is to detect and segment the hand from the background, which can be difficult in the wild due to the varying backgrounds, skin colours, and lighting conditions Shi et al. (2019). This is a challenge for models that are not robust to these conditions, and can lead to poor performance.

### 3.7. *Overcoming Obstacles (RQ-5)*

- Propose techniques to enhance accuracy and address the challenges identified in the previous section.
- Discuss data augmentation strategies and the use of transfer learning and pre-trained models.

### 3.8. *State of the Art and Real-World Applications (RQ-4, RQ-7)*

Automatic Speech Recognition (ASR) is a separate domain from Automatic Sign Language Recognition (ASLR). ASR is the task of converting speech to text and is a much more mature field. While one is an audio task, and the other a visual one, they both involve converting a sequence of symbols to text. As there is more research throughput into ASR, it's very common to see overlapping models and models adapted from ASR to ASLR. An example of this is the Conformer-CTC model Gulati et al. (2020) and the more recent advancement Squeezeformer Kim et al. (2022). Although

academia is advancing the field of ASLR formally, there are many active practitioners experimenting on websites through Kaggle competitions Manfred Georg, Mark Sherwood, Phil Culliton, Sam Sepah, Sohier Dane, Thad Starner Ashley Chow, Glenn Cameron (2023). Whilst this is not a formal academic setting, it is a good indicator of the state of the art, and the models that are being used in the real world. Furthermore, it mustn't be understated that these models and results aren't peer-reviewed, and are not necessarily reproducible. However, these developments are on the internet immediately, and aren't being kept back as research papers that could take months and years to publish. They could be used to inform the development of our own model, given the time constraints of this project.

\*\*\*\*\* STATE OF THE ART \*\*\*\*\* Transformers/Squeezeformer Beam Search with CTC decoder. Multi-headed Attention Mechanisms CTC loss. An insane amount of microadjustments/augmentations to the data to optimize training and score.

- Present the most recent and influential works in ASL fingerspelling recognition.
- Highlight real-world applications, with a focus on practical aspects like kiosk systems.

### 3.9. Future Directions and Open Challenges (RQ-7)

- Discuss emerging trends in the field of ASL fingerspelling recognition.
- Identify areas that require further research and exploration.
- Explore the potential impact of future advancements, particularly in deep learning.

### 3.10. Ethical and Societal Considerations

- Address ethical considerations, including data privacy concerns.
- Discuss issues related to bias and fairness in ASL recognition models.
- Examine the implications of ASL fingerspelling recognition technology for the deaf and hard of hearing community.

### 3.11. Conclusion

- Summarize the main findings related to each research question discussed throughout the report.
- Reiterate the importance of the topic and its potential impact on the field of ASL fingerspelling recognition.

## References

- Abadi, M., Barham, P., Chen, J., Chen, Z., Davis, A., Dean, J., Devin, M., Ghemawat, S., Irving, G., Isard, M., Kudlur, M., Levenberg, J., Monga, R., Moore, S., Murray, D.G., Steiner, B., Tucker, P., Vasudevan, V., Warden, P., Wicke, M., Yu, Y., Zheng, X., 2016. TensorFlow: A system for large-scale machine learning. doi:10.48550/arXiv.1605.08695, arXiv:1605.08695.
- Abiyev, R., Idoko, J.B., Arslan, M., 2020. Reconstruction of Convolutional Neural Network for Sign Language Recognition, in: 2020 International Conference on Electrical, Communication, and Computer Engineering (ICECCE), IEEE, Istanbul, Turkey. pp. 1–5. doi:10.1109/ICECCE49384.2020.9179356.
- Bahdanau, D., Cho, K., Bengio, Y., 2016. Neural Machine Translation by Jointly Learning to Align and Translate. doi:10.48550/arXiv.1409.0473, arXiv:1409.0473.
- Bantupalli, K., Xie, Y., 2018. American Sign Language Recognition using Deep Learning and Computer Vision, in: 2018 IEEE International Conference on Big Data (Big Data), pp. 4896–4899. doi:10.1109/BigData.2018.8622141.
- Chong, T.W., Lee, B.G., 2018. American Sign Language Recognition Using Leap Motion Controller with Machine Learning Approach. Sensors 18, 3554. doi:10.3390/s18103554.
- Cihan Camgoz, N., Koller, O., Hadfield, S., Bowden, R., 2020. Sign Language Transformers: Joint End-to-End Sign Language Recognition and Translation, in: 2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), IEEE, Seattle, WA, USA. pp. 10020–10030. doi:10.1109/CVPR42600.2020.01004.
- Culjak, I., Abram, D., Pribanic, T., Dzapo, H., Cifrek, M., . A brief introduction to OpenCV .
- Ethnologue, 2023. American Sign Language — Ethnologue Free. URL: <https://www.ethnologue.com/language/ase/>.
- Gajurel, K., Zhong, C., Wang, G., 2021. A Fine-Grained Visual Attention Approach for Fingerspelling Recognition in the Wild, in: 2021 IEEE International Conference on Systems, Man, and Cybernetics (SMC), IEEE, Melbourne, Australia. pp. 3266–3271. doi:10.1109/SMC52423.2021.9658982.
- Graves, A., Fernández, S., Gomez, F., Schmidhuber, J., 2006. Connectionist temporal classification: Labelling unsegmented sequence data with recurrent neural networks, in: Proceedings of the 23rd

- International Conference on Machine Learning, Association for Computing Machinery, New York, NY, USA. pp. 369–376. doi:10.1145/1143844.1143891.
- Gulati, A., Qin, J., Chiu, C.C., Parmar, N., Zhang, Y., Yu, J., Han, W., Wang, S., Zhang, Z., Wu, Y., Pang, R., 2020. Conformer: Convolution-augmented Transformer for Speech Recognition. doi:10.48550/arXiv.2005.08100, arXiv:2005.08100.
- Kabade, A.E., Desai, P., C, S., G, S., 2023. American Sign Language Fingerspelling Recognition using Attention Model, in: 2023 IEEE 8th International Conference for Convergence in Technology (I2CT), pp. 1–6. doi:10.1109/I2CT57861.2023.10126277.
- Kim, S., Gholami, A., Shaw, A., Lee, N., Mangalam, K., Malik, J., Mahoney, M.W., Keutzer, K., 2022. Squeezeformer: An Efficient Transformer for Automatic Speech Recognition. doi:10.48550/arXiv.2206.00888, arXiv:2206.00888.
- Krizhevsky, A., Sutskever, I., Hinton, G.E., 2012. ImageNet Classification with Deep Convolutional Neural Networks, in: Advances in Neural Information Processing Systems, Curran Associates, Inc.. pp. 1097–1105. URL: [https://proceedings.neurips.cc/paper\\_files/paper/2012/hash/c399862d3b9d6b76c8436e924a6](https://proceedings.neurips.cc/paper_files/paper/2012/hash/c399862d3b9d6b76c8436e924a6)
- Lecun, Y., Bottou, L., Bengio, Y., Haffner, P., 1998. Gradient-based learning applied to document recognition. Proceedings of the IEEE 86, 2278–2324. doi:10.1109/5.726791.
- Lugaresi, C., Tang, J., Nash, H., McClanahan, C., Uboweja, E., Hays, M., Zhang, F., Chang, C.L., Yong, M.G., Lee, J., Chang, W.T., Hua, W., Georg, M., Grundmann, M., 2019. MediaPipe: A Framework for Building Perception Pipelines. doi:10.48550/arXiv.1906.08172, arXiv:1906.08172.
- Manfred Georg, Mark Sherwood, Phil Culliton, Sam Sepah, Sohier Dane, Thad Starner Ashley Chow, Glenn Cameron, 2023. Google - american sign language fingerspelling recognition. URL: <https://kaggle.com/competitions/asl-fingerspelling>.
- Mitchell, R.E., Young, T.A., Bachelda, B., Karchmer, M.A., 2006. How Many People Use ASL in the United States?: Why Estimates Need Updating. Sign Language Studies 6, 306–335. URL: <https://www.jstor.org/stable/26190621>, arXiv:26190621.
- Munib, Q., Habeeb, M., Takruri, B., Al-Malik, H.A., 2007. American sign language (ASL) recognition based on Hough transform and neural networks. Expert Systems with Applications 32, 24–37. doi:10.1016/j.eswa.2005.11.018.

- Nguyen, H.B., Do, H.N., 2019. Deep Learning for American Sign Language Fingerspelling Recognition System, in: 2019 26th International Conference on Telecommunications (ICT), IEEE, Hanoi, Vietnam. pp. 314–318. doi:10.1109/ICT.2019.8798856.
- Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z., Gimelshein, N., Antiga, L., Desmaison, A., Köpf, A., Yang, E., DeVito, Z., Raison, M., Tejani, A., Chilamkurthy, S., Steiner, B., Fang, L., Bai, J., Chintala, S., 2019. PyTorch: An Imperative Style, High-Performance Deep Learning Library. doi:10.48550/arXiv.1912.01703, arXiv:1912.01703.
- S Kumar, S., Wangyal, T., Saboo, V., Srinath, R., 2018. Time Series Neural Networks for Real Time Sign Language Translation, in: 2018 17th IEEE International Conference on Machine Learning and Applications (ICMLA), IEEE, Orlando, FL. pp. 243–248. doi:10.1109/ICMLA.2018.00043.
- Saeed, Z.R., Zainol, Z.B., Zaidan, B.B., Alamoodi, A.H., 2022. A Systematic Review on Systems-Based Sensory Gloves for Sign Language Pattern Recognition: An Update From 2017 to 2022. IEEE Access 10, 123358–123377. doi:10.1109/ACCESS.2022.3219430.
- Sherstinsky, A., 2020. Fundamentals of Recurrent Neural Network (RNN) and Long Short-Term Memory (LSTM) Network. Physica D: Nonlinear Phenomena 404, 132306. doi:10.1016/j.physd.2019.132306, arXiv:1808.03314.
- Shi, B., Brentari, D., Shakhnarovich, G., Livescu, K., 2021. Fingerspelling Detection in American Sign Language, in: 2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), IEEE, Nashville, TN, USA. pp. 4164–4173. doi:10.1109/CVPR46437.2021.00415.
- Shi, B., Brentari, D., Shakhnarovich, G., Livescu, K., 2022. Searching for fingerspelled content in American Sign Language. URL: <http://arxiv.org/abs/2203.13291>, arXiv:2203.13291.
- Shi, B., Del Rio, A.M., Keane, J., Michaux, J., Brentari, D., Shakhnarovich, G., Livescu, K., 2018. American Sign Language Fingerspelling Recognition in the Wild, in: 2018 IEEE Spoken Language Technology Workshop (SLT), pp. 145–152. doi:10.1109/SLT.2018.8639639.
- Shi, B., Rio, A.M.D., Keane, J., Brentari, D., Shakhnarovich, G., Livescu, K., 2019. Fingerspelling Recognition in the Wild With Iterative Visual Attention, in: 2019 IEEE/CVF International Conference on Computer Vision (ICCV), IEEE, Seoul, Korea (South). pp. 5399–5408. doi:10.1109/ICCV.2019.00550.
- Srinivasan, R., Kavita, R., Kavitha, M., Mallikarjuna, B., Bhatia, S., Agarwal, B., Ahlawat, V., Goel, A., 2023. Python And Opencv For Sign Language Recognition, in: 2023 International Confer-

ence on Device Intelligence, Computing and Communication Technologies, (DICCT), pp. 1–5. doi:10.1109/DICCT56244.2023.10110225.

Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, L., Polosukhin, I., 2023. Attention Is All You Need. doi:10.48550/arXiv.1706.03762, arXiv:1706.03762.

Vogler, C., Metaxas, D., 1999. Parallel hidden Markov models for American sign language recognition, in: Proceedings of the Seventh IEEE International Conference on Computer Vision, pp. 116–122 vol.1. doi:10.1109/ICCV.1999.791206.

von Agris, U., Zieren, J., Canzler, U., Bauer, B., Kraiss, K.F., 2008. Recent developments in visual sign language recognition. Universal Access in the Information Society 6, 323–362. doi:10.1007/s10209-007-0104-x.

Weerasooriya, A.A., Ambegoda, T.D., 2022. Sinhala Fingerspelling Sign Language Recognition with Computer Vision, in: 2022 Moratuwa Engineering Research Conference (MERCon), IEEE, Moratuwa, Sri Lanka. pp. 1–6. doi:10.1109/MERCon55799.2022.9906281.

## 4. MAIN CHAPTERS:

### 4.1. Chapter 1:

## 5. DEVELOPMENT NOTES

### 5.1. Algorithm Pseudocode

<https://ctan.math.washington.edu/tex-archive/macros/latex/contrib/algpseudocodex/algpseudocodex.pdf>

### 5.2. Model Art

<https://github.com/ashishpatel26/Tools-to-Design-or-Visualize-Architecture-of-Neural-Network?tab=readme-ov-file>

Appendices

**Appendix A. Glossary**



**Appendix B. Marking Scheme**

## **Appendix C. Changes to the Project Initiation Document**

**Appendix D. Current Environment Investigation Report**

## Appendix E. Requirements Specification

## Appendix F. Design Report

## Appendix G. Implementation

## Appendix H. Testing

## **Appendix I. User Guide**



## Appendix J. Project Management

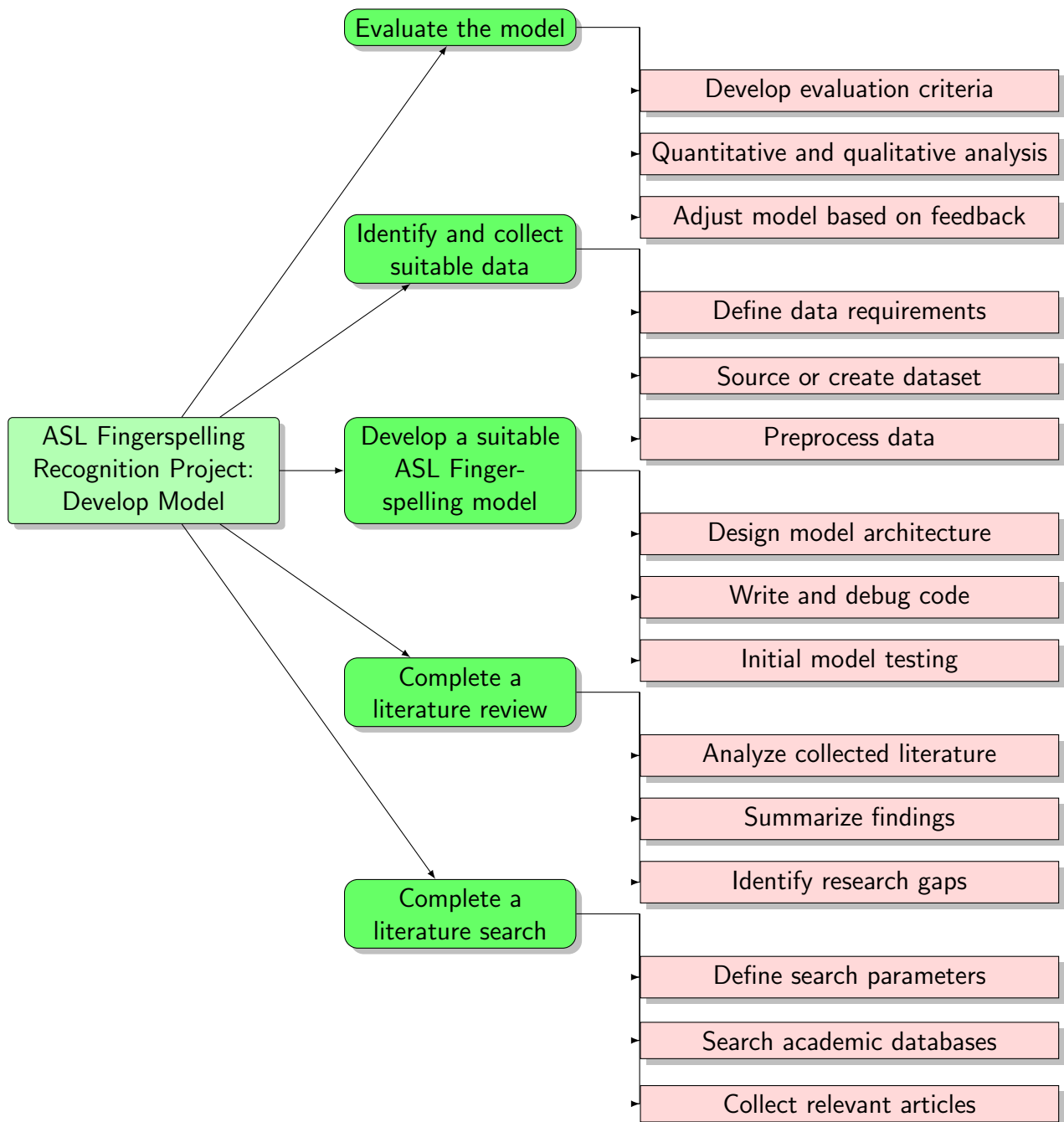


Figure J.1: Work Breakdown Structure: Develop Model

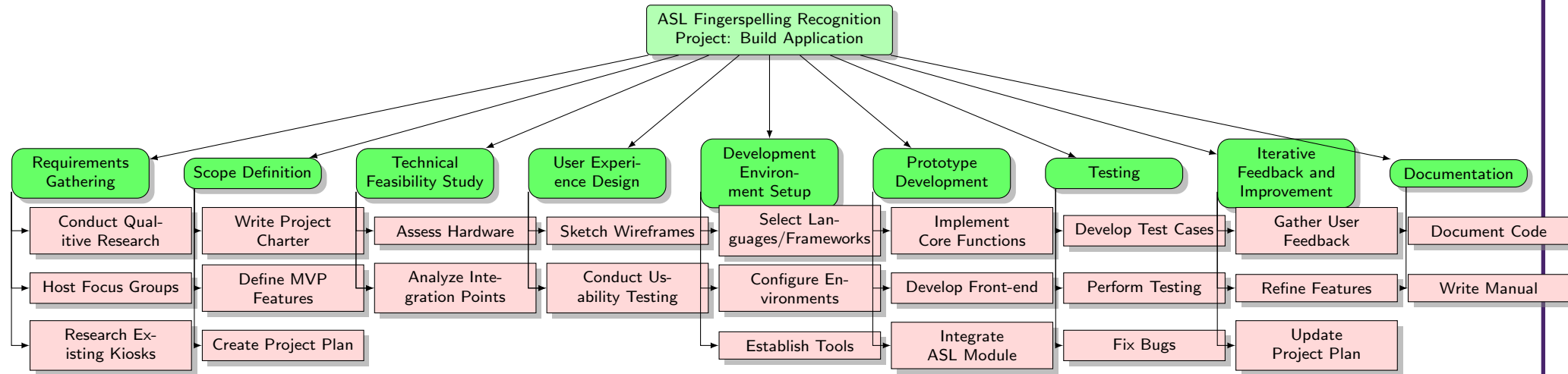


Figure J.2: Work Breakdown Structure: Build Application

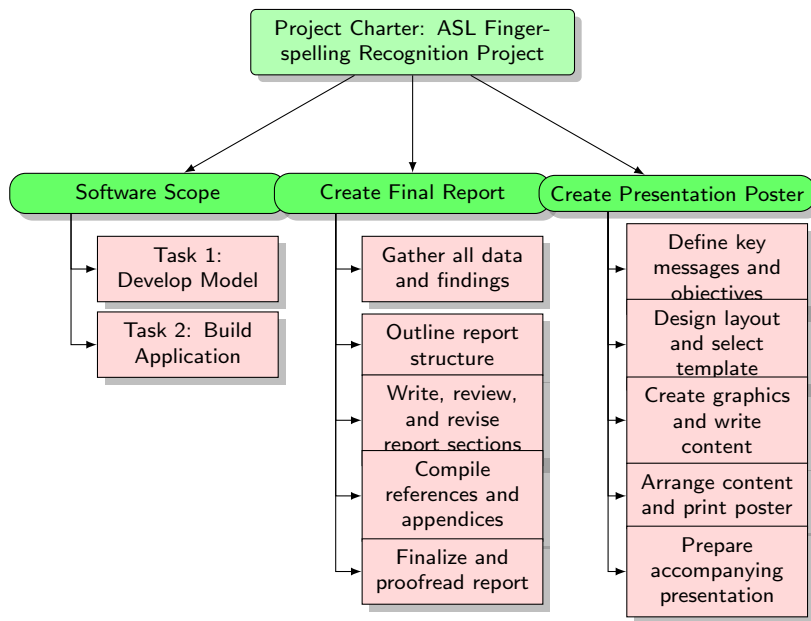


Figure J.3: Work Breakdown Structure

Gantt chart placeholder  
(Refer to the accompanying zipped file for the full chart)

Figure J.4: Gantt Chart

In this report's accompanying zipped file, a detailed Gantt chart is included as a separate document. Due to its extensive size and complexity, it is provided as an individual file to facilitate detailed review and ensure clarity. Please refer to the zipped file named gantt-project.png for the complete Gantt chart, which offers an in-depth view of the project timeline and milestones.

<b>Risk</b>	<b>Impact</b>	<b>Probability</b>	<b>Status</b>	<b>Mitigation Strategy</b>
Inaccurate ASL Recognition	High	Medium	Open	Enhance data collection and improve algorithm accuracy.
Data Privacy Concerns	High	High	Open	Implement GDPR compliant data handling processes.
Loss of Data	High	Low	Open	Implement GoogleDrive and GitHub repository.
Loss of Project Supervisor	High	Low	Open	Maintain regular communication with supervisor.
Project Delays	Medium	High	Open	Develop a schedule with buffers and regularly update it.
User Adoption Challenges	Medium	High	Open	Engage with users early and incorporate feedback.
Technology Integration Issues	Medium	Medium	Open	Conduct compatibility testing.

Figure J.5: Risk register

**Appendix K. Meetings With Supervisor**

<b>Date</b>	<b>Time</b>	<b>Location</b>	<b>Purpose</b>	<b>Description and Actions</b>
15 November 2023	13:00-14:00	CCCU - Lg33	Discuss project and literature research	Reviewed current status, discussed challenges, and agreed on next steps including further research on ASL recognition and user-centric approach to requirements. Contact charities
29 November 2023	13:00-14:00	CCCU - Lg33	Discuss third party inclusion, user centric requirement	Update on contacted charities and double diamond requirements
13 December 2023	13:00-13:30	Online	Catchup and talk poster presentation	Spoke about markschemes, poster style, sharing work

## **Appendix L. Agile Development: Timebox 1**

**Appendix M. Agile Development: Timebox 2**



**Appendix N. Agile Development: Timebox 3**