# Neuro-symbolic Answer Set Programming for Human Activity Recognition in Videos

Andreas Oikonomakis

Athens, September 2023

# Acknowledgements

## Περίληψη

Η μηχανική μάθηση (Machine / Deep learning) και ο συλλογισμός μέσω ανα-παράστασης γνώσης (Machine reasoning) θεωρούνται δυο διαφορετικά πεδία της Τεχνητής Νοημοσύνης. Μέσω της μηχανικής μάθησης επιτυγχάνουμε τη δημιουρ-γία μοντέλων που έχουν αντιληπτικές δυνατότητες χαμηλού επιπέδου όπως ταξι-νόμηση και πρόβλεψη ενώ μέσω του συλλογισμού με τη χρήση κατάλληλων ανα-παραστάσεων και τεχνικών λογικής (logic based reasoning) επιτυγχάνουμε την εξαγωγή συμπερασμάτων και πληροφορίων σε υψηλότερο επίπεδο. Ο συνδυασμός αυτών των δυο διαφορετικών πεδίων (neural-based learning , logic-based reason-ing) θα μπορούσε να συμβάλει στη δημιουργία νέων συστημάτων τα οποία θα είναι ικανά να έχουν αντίληψη του περιβάλλοντος και να μπορούν να εξάγουν συμπε-ράσματα πάνω στα δεδομένα που τους έχουν δοθεί. Στην παρούσα διπλωματική θα επικεντρωθούμε στο Neural-symbolic computation με το οποίο επιτυγχάνεται ο συνδυασμός deep-learning και συλλογισμού μέσω μιας υπάρχουσας μεθόδου και λογισμικού με το όνομα NeurAsp. Αρχικά μέσω απλών παραδειγμάτων θα παρου-σιάσουμε τις δυνατότητες αυτής της μεθόδου και πως αυτή λειτουργεί εσωτερικά και ενσωματώνεται με τις παραδοσιακές μεθόδους μηχανικής μάθησης και στην συνέχεια θα εφαρμόσουμε την μέθοδο αυτή στην αναγνώριση δραστηριοτήτων σε βίντεο μεταξύ ανθρώπων χρησιμοποιώντας παράλληλα τεχνικές αναγνώρησης πε-ρίπλοκων γεγονότων Complex Event Recognition (CER). Συγκεκριμένα μέσω τριών πειραμάτων στα οποία συγκρίνουμε τόσο τις παραδοσιακές μεθόδους βαθιάς μηχανικής μάθησης όσο και αυτές του NeurAsp θα επιχειρήσουμε να αναδείξουμε τα ωφέλη των μεθόδων που συνδυάζουν μηχανισμούς λογικής μαζί με την βαθιά μηχανική μάθηση.

# Abstract

Machine / Deep Learning and Machine Reasoning are considered two different subfields of Artificial Intelligence. With machine learning methods we can build models with low level perceptual capabilities and with logic based methods we can extract information and perform reasoning at a higher level. Combining neural learning methods with logic-based techniques could help create systems that are able to perceive their environment and infer the data given as input. In this thesis, we will focus on neurosymbolic computation, where the combination of deep learning and reasoning is achieved through an existing framework called NeurAsp. We will go through simple examples that demonstrate NeurAsp's capabilities and show how it works and integrates internally with traditional deep learning methods. The main goal of this thesis is to apply this method to the task of detecting human activity in videos with the usage of Complex Event Recognition (CER) techniques. Finally, we will show the benefits of integrating logic-based techniques with neural methods by presenting three different experimental setups in which we compare the performances of pure traditional deep learning methods and those proposed by the NeurAsp framework.

# Contents

# 1 Introduction

## 1.1 Neuro-symbolic Artificial Intelligence

Neuro-symbolic artificial intelligence refers to a field of research that combines neural and symbolic methods, which are considered distinct areas of artificial intelligence. Neural methods usually involve training artificial neural networks using large amounts of raw data, where the model learns to associate the given data by continuously learning from its mistakes and eventually is able to make predictions for the given task even for new data related to the task. For example if we train a neural network to identify a set of animals from images such as dogs, cats and birds we can give as input a new instance of a bird and get the results with high confidence. This fact shows that neural nets are very capable at capturing patterns from raw data by adjusting the strengths of the connections between its nodes and finally be able to infer on a new image. The strength of neural nets depends heavily on the amount of training data and acquiring training data is often costly and sometimes impossible. Also it is known that neural nets are sensitive to noise even if the noise is obscure to the human eye. On the other hand, symbolic methods involve the explicit representation of knowledge using formal languages, logic-based techniques and reasoning allow us to infer and extract information at a higher level. For instance if we encode information about the color, sizes and shapes of objects to a symbolic program and create a general rule regarding the similarities such as "the objects are similar if they have the same size or color" we could answer questions such as "are the square and the triangle similar?" given that squares and triangles are encoded as symbols. Symbolic AI employs an inference engine, which uses rules of logic to answer queries but this requires a lot of effort for encoding this kind of information and most important it requires domain expertise which denotes the major downside of symbolic AI systems. The symbolic AI system cannot learn on its own: for instance if we have encoded only squares and triangles in

our example and we ask about circles the system will fail because circles are not present in the knowledge base. Given the strengths and weaknesses of those two different AI methods, neuro-symbolic AI is a novel area of AI research which seeks to create hybrid methods with promising results.[Susskin, 2021]. Logic-based methods are able to encode domain knowledge and also provide proof regarding the correctness of the result allowing us to have better explainability which will benefit neural methods to reduce their black-box nature. Also with logic-based methods we can provide some sort of guidance during the training phase of a neural network by integrating logic rules and not relying only to the traditional neural loss functions leading to less time consuming and data demanding training procedure. On the other hand neural methods with their ability to handle and model large amount of raw data could be used to efficiently approximate combinatorial optimization tasks which is the major drawback of logic-based methods. With those hybrid approaches each method can benefit from the other since neuro-symbolic AI aims to combine the complementary strengths of neural and logic based methods.

## 1.2   Neuro-symbolic Approaches

There are several discussions and proposals on how to integrate neural and logic methods which can be categorized to set of neuro-symbolic (NeSy) families. One of the earliest NeSy approaches was to train neural networks to mimic the behavior of logical reasoners(Logic as Neural Programs). The neurons of a neural network can operate as logical gates(AND, OR ,etc..). This way neurons can represent logical connectives, a piece of a logical statement or a concept.[Rocktäschel, 2017][Sourek, 2018][Evans, 2018] Those approaches provided an alternative way of logical reasoning by tackling some of the weaknesses of logic-based methods such as high computational complexity, proneness to minor inconsistencies in data and finally the dependency on domain expertise but the black-box nature of neural networks nullify some of the benefits of logic-based methods such as their expressive power. Another NeSy approach is to train a neural network to comply with a set of logical constraints(Logic as regularizer) by incorporating logical constraints in the loss function that penalizes violations of logical rules.[Diligenti, 2017][Xu, 2018][Marra, 2021][Badreddine, 2022]. The logic as regularizer NeSy approaches often introduce additional layers to perform logic operations keeping the logical and neural parts separated leading to Nesy Interfacing approaches.[Manhaeve, 2018][Yang, 2020][Winters, 2022]. In NeSy interfacing a logic reasoner takes as input the predictions of a neural network, computes the gradients in respect of the defined logic rules and then communicates the gradient updates via the back-propagation which is the usual training procedure of a neural network. This clear separation between the neural and logic parts is very useful when we want to interpret the behaviour of a neural network but it adds an overhead to training and inference times.

In Neural Probabilistic Logic Programming [Manhaeve, 2018] the authors treat outputs of a neural network as the probability distribution over atomic facts and based on this idea NeurAsp is proposed. NeurAsp [Yang, 2020] is a

Figure 1: NeSy Approaches

framework [1] that focuses on Neural-symbolic computation [Garcez, 2019] which achieves the combination of deep-learning and reasoning based on Answer Set Programming [Lifschitz, 2019]. Answer Set Programming (ASP) is a method of solving combinatorial optimization problems, based on logic programming which is supported by specialized tools (ASP solvers). ASP solvers allow the combination of logic based reasoning with optimization methods and also various tasks related to artificial intelligence such as reasoning under uncertainty and symbolic learning can be expressed and solved by using ASP solvers.

## 1.3 Human Activity Recognition (HAR)

Human Activity Recognition (HAR) is a challenging time series classification task that involves predicting a person's movement based on input data that comes from various devices such as gps, wearable sensors, video cameras etc. HAR is a multi staged task [Gupta, 2022] starting from capturing the signal or video from the device, perform any data pre-processing needed and finally choose the suitable machine or deep learning algorithm to train. One of the most studied applications of HAR is the video-based HAR which aims to automatically recognize and categorize human activities in video sequences. There are numerous deep learning-based methods for activity recognition in videos that mainly involve Convolutional Neural Networks (CNNs) [Dogan, 2021] and Recurrent Neural Networks (RNNs) [Sungh, 2018] or their combination the three-dimensional CNNs (C3Ds)[Wang, 2015]. Each method has its benefits and limitations due to the complex nature of this particular task since the de-

---

[1]https://github.com/azreasoners/NeurAsp

velopers of an automatic HAR system have to face a lot of problems such as background blur, partial occlusion, image quality resolution, changes in scale and viewpoint and so on. Finally, the annotation of behavioral roles is time-consuming and costly, resulting in the unavailability of HAR datasets, and the similarities between action classes also result in significant amount of noise in the available datasets.[Vrigkas, 2015]

## 1.4 Complex Event Recognition in HAR

Human activity recognition as mentioned involves various time series data sources such as sensors , cameras etc. If we treat the input as a timestamped information and put a meaningful label we can construct entities called simple events. Those simple events called primitives is assumed that they are not dependent on other events. Some events could be the result of an external pre-process procedure such as the average speed of a person in a given time interval. Those events could contain patterns that are useful to extract further insights on our data. Complex Event Recognition (CER) techniques allow the identification of meaningful patterns in temporal data by creating complex event definitions that may contain occurrences of simple events, combinations of simple events or even combinations of other complex events.[Katzouris, 2017] In HAR it is important to model the occurrences, the duration and the effects of the events. For this reason the usage of a logic-based CER system with the ability to reason on temporal data is essential.[Cugola, 2012] [Giatrakos, 2019] For example lets assume that two people are walking in a very close distance and a certain orientation for a time interval, then potentially a complex event "walking together" is triggered and after some time one of them starts running away. There must be a logic mechanism that would trigger the termination of the "walking together" by the time one of those two persons starts running. The Event Calculus dialect as proposed in RTEC [Artikis, 2015] offers a solution to perform temporal reasoning and build a suitable logic-based CER system for HAR.[Alevizos, 2022]

## 2 Asp And Event Calculus Background

At this point we will provide the needed background for answer set programming , the event calculus and how they can be used in the CER context as provided by [Katzouris, 2023].

**Answer Set Programming**: In what follows a rule $r$ is an expression of the form $\alpha \leftarrow \delta_1, ..., \delta_n$, where $\alpha$ is an atom, called the head of $r$, $\delta_i$s are literals (possibly negated atoms), which collectively form the body of $r$ and commas in the bodies of rules denote conjunction. A rule is ground if it contains no variables and a grounding of a rule $r$ is called an instance of $r$. A (Herbrand) interpretation is a collection of true ground facts. An interpretation $I$ satisfies an atom $\alpha$ iff $\alpha \in I$. $I$ satisfies a ground rule iff satisfying each literal in the body implies that the head atom is also satisfied and it satisfies a non-ground rule $r$ if it satisfies all ground instances of r. An interpretation $I$ is a model of a logic

program $\pi$ (collection of rules) if it satisfies every rule in $\pi$ and it is a minimal model if no strict subset of I has this property. An interpretation $I$ is an answer set of $\pi$ if it is a minimal model of the reduct of $\pi$, that is, the negation-free, ground program that results by removing from the ground version of $\pi$ all rules with a negated body literal not satisfied by $I$ and removing all negated literals from the bodies of the remaining rules

**The Event Calculus** is a temporal logic for reasoning about events and their effects. Its ontology comprises time points (integers), fluents, that is, properties which have certain values in time, and events, that is, occurrences in time that may affect fluents and alter their value. Its axioms incorporate the commonsense law of inertia, according to which fluents persist over time, unless they are affected by an event. Its basic predicates and axioms are presented in Table 1(a), (b). Axiom (1) states that a fluent $F$ holds at time $T$ if it has been initiated at the previous time point, while Axiom (2) states that $F$ continues to hold unless it is terminated. Definitions of *initiatedAt/2* and *terminatedAt/2* predicates are provided in an application-specific manner. Using the Event Calculus in a CER context allows to reason with CEs that have duration in time and are subject to commonsense phenomena, via associating CEs to fluents. In this case, a set of CE patterns is a set of *initiatedAt/2* and *terminatedAt/2* rules. As an example we use the task of activity recognition, as defined in the CAVIAR project[2]. The CAVIAR data set consists of videos of a public space, where actors per Table 1. (a), (b) The basic predicates and the Event Calculus axioms. (c) Example CAVIAR data. At time point 1 person with id1 is walking, her $(X, Y)$ coordinates are (201, 454) and her direction is 270°. The target CE atoms (true state – supervision) for time point 1 state that persons id1 and id2 are moving together at the next time point. (d) An example of two domain-specific axioms in the EC. For example, the first rule dictates that moving together between two persons $X$ and $Y$ is initiated at time $T$ if both $X$ and $Y$ are walking at time $T$, their euclidean distance is less than 25 pixel positions and their difference in direction is less than 45°. The second rule dictates that moving together between $X$ and $Y$ is terminated at time $T$ if one of them is standing still at time $T$ and their euclidean distance at $T$ is greater than 30 form some activities. These videos have been manually annotated by the CAVIAR team to provide the ground truth for two types of activity. The first type, corresponding to simple events, consists of knowledge about a person's activities at a certain video frame/time point (e.g. walking, standing still and so on). The second type, corresponding to CEs/fluents, consists of activities that involve more than one person, for instance two people moving together, meeting each other and so on. The aim is to detect CEs as combinations of simple events and additional domain knowledge, such as a person's position and direction. Table 1(c) presents an example of CAVIAR data, consisting of observations for a particular time point, in the form of an interpretation $I_1$. A stream of interpretations is matched against a set of CE patterns (initiation/termination rules – see Table 1(d)), to infer the truth values of CE

---

[2]https://homepages.inf.ed.ac.uk/rbf/CAVIARDATA1

instances in time, using the Event Calculus axioms as a reasoning engine. We henceforth call the atoms corresponding to CE instances whose truth values are to be inferred/predicted, target CE instances. Table 1(c) presents the target CE instances corresponding to the observations in $I_1$.

**(a)**

| Predicate: | Meaning: |
|---|---|
| *happensAt(E,T)* | Event $E$ occurs at time $T$ |
| *initiatedAt(F,T)* | At time $T$, a period of time for which fluent $F$ holds is initiated. |
| *terminatedAt(F,T)* | At time $T$, a period of time for which fluent $F$ holds is terminated. |
| *holdsAt(F,T)* | Fluent $F$ holds at time $T$ |

**(b) Axioms of the Event Calculus**

$$holdsAt(F,T+1) \leftarrow initiatedAt(F,T). \tag{1}$$

$$holdsAt(F,T+1) \leftarrow holdsAt(F,T), not\ terminatedAt(F,T). \tag{2}$$

**(c)**

| Observations $I_1$ at time 1: | Target CE instances at time 1: |
|---|---|
| {*happensAt(walk(id_1,1)*, | {*holdsAt(move(id_1,id_2),2)*, |
| *happensAt(walk(id_2,1)*, | *holdsAt(move(id_2,id_1),2)* } |
| *coords(walk(id_1,201,454,1)*, | |
| *coords(walk(id_2,230,440,1)*, | |
| *coords(direction(id_1,270,1)*, | |
| *coords(direction(id_2,270,1)*, } | |

<div align="center">Table 1</div>

CER techniques require simple events in order to operate and those simple events must be extracted from collected raw data. In HAR for instance raw data could be the coordinates of the bounding boxes of persons provided by a neural network that is trained to detect persons but there is a big gap between the raw coordinates and a concept such as "walking" or "running".This fact indicates

the importance of neuro-symbolic techniques since we can bridge this gap by using deep learning techniques where we can train a neural network to identify and extract those simple events from raw data and provide them as inputs to a logic-based CER system.

# 3 Basic Deep Learning Concepts

In this section we will introduce some basic concepts that are used in the training procedure of a neural model. Of course it is not an exhaustive list of definitions but rather a simple reference to some terms that we will use in the following sections of this thesis. We will focus on the supervised training algorithm where given an input and an output the neural network learns to associate them.

## 3.1 Neural Network Layers and Neurons

A neural network can be described by three types of layers :

- **Input Layer** : This layer will accept the raw data and pass it to the rest of the network

- **Hidden Layer** : This type of layer may have more than one instance in a neural network and can perform multiple mathematical operations and data transformations.

- **Output Layer** : The final layer in the neural network where predictions are obtained.



Figure 2: NN Visual Representation

A layer consists of small individual units called **neurons**. Neurons can take an arbitrary number of inputs and have an output given condition. A simple

example is given below where a neuron takes as input two values X1, X2 and outputs a value Y.



Figure 3: Neuron Visual Representation

When a neuron provides an output it is common to use the phrase **"the neuron fired"**. A neuron fires when a condition is met. For instance this condition in our depicted example could be a simple expression such as the "sum of inputs must be greater than 10 in order to fire". We can assume when the neuron fires the output Y is 1 and when it doesn't is 0:



Figure 4: Neuron operations given a condition

Here we will introduce the terms **bias** and **weights**. In our example the threshold was 10 and it is a constant that we introduced to the neuron in order to create a condition for the inputs.

$$X_1 + X_2 > threshold \rightarrow \textcolor{green}{Y=1}$$
$$X_1 + X_2 < threshold \rightarrow \textcolor{red}{Y=0}$$

if we move the threshold to the other side we get:

$$X_1 + X_2 - threshold > 0$$
$$X_1 + X_2 - threshold < 0$$

Now we can rename the negative threshold with the positive term bias:

$$X_1 + X_2 + bias > 0$$
$$X_1 + X_2 + bias < 0$$

So far our inputs have the same "importance" and impact to the outcome of the neuron but if we multiply each term with a number called **weight** we can quantify the impact of each input:

$$W_1 * X_1 + W_2 * X_2 + W_0 * bias > 0$$

## 3.2    Activation Functions

A neuron can output any number in the range of $(-\infty \; \infty)$ We use activation functions to limit the range of the output of a neuron. This helps with the internal computations between the layers of a neural network and also introduces non linearity which is essential for the learning procedure. We introduce some of the most well known activation functions:

- **Sigmoid activation function** : The sigmoid function has its domain as the set of all real numbers, and its range is (0, 1). In other words this function has an output close to 0 for large negative input values and an input close to 1 for large positive values.



Figure 5: Sigmoid activation function

- **Relu activation function** : The relu function is a simple calculation that returns the value provided as input directly, or the value 0 if the input is 0 or less. This function is the most used activation function between the hidden layers of a neural network.



$$\text{ReLU}(z) = \begin{cases} z, z > 0 \\ 0, otherwise \end{cases}$$

Figure 6: Relu activation function

- **Tanh activation function** : This function has its domain as the set of all real numbers, and its range is (-1, 1). A main difference with the sigmoid activation function is that this function introduces negative numbers to its output.Its output is a value close to -1 for large negative input values and an input close to 1 for large positive values.



$$\sigma(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

Figure 7: Tanh activation function

- **Softmax activation function** : Activation functions are also used in the output layer of the neural network in order to produce the desired outcome depending on the task. The softmax activation function transforms a previous layer's output into a vector of probabilities. It is commonly used for multiclass classification.

Figure 8: Softmax activation function

## 3.3 Loss Functions

The results of the Output Layer are compared to the ground truth targets provided by the dataset annotators.The loss function is a quantitative indicator of the performance of a neural network that computes the difference between the prediction and the actual target value. The difference between the prediction $\hat{Y}$ and the actual value Y is called **loss** or **error**. One of the most common loss function is the **Mean Square Error (MSE)**:

$$MSE = \frac{1}{n} \sum_{i=1}^{n} (\hat{Y}_i - Y_i)^2$$

The MSE loss function penalizes the model for making large errors by squaring them. This loss function is used in regression tasks where the output needs to be a numerical value
For classification Tasks a very common loss function is the **Cross-entropy loss (Log Loss)**:

$$L = -\frac{1}{n} \sum_{i=1}^{n} Y_i \cdot log(\hat{Y}_i)$$

We use this loss function to measure the performance of a classification model whose output is a probability value between 0 and 1.

## 3.4 Internal Training Steps of a Neural Network

The training procedure of a neural network is a set of steps that are repeated many times. Here we provide a brief description of those steps:

- **Forward Pass**: During the forward pass the raw input data in the input layer traverses through all the hidden layers where a set of mathematical operations and transformations is being performed and finally the result reaches the output layer as the prediction of the neural network

- **Loss computation**: The loss function computes the loss between the prediction and the ground truth.

- **Gradient computation and weight updates**: The **gradient descent algorithm** is used to find the local minimum of a function and the **gradient ascent** is used to find the local maximum of a function. In any of those algorithms the gradients of a function are computed and multiplied by a given step called the **learning rate**. The values that come from this multiplication are the needed updates for the weights.

- **Backward Pass (Backpropagation)**:The backward pass is most essential step of the learning procedure since it is responsible to apply the weight updates to the hidden layers.

  As mentioned those steps are repeated for all samples in the training dataset and the process of all samples is called an **epoch**. Due to computer memory limitations is if often that the data is splitted to smaller chunks called **bathes**. A neural network can train for several epochs until some criteria are met. (for instance the loss cannot be minimized any further)

## 3.5  Performance Metrics for Classification Tasks

During or after the training procedure it is essential to evaluate the performance of the neural network.The most intuitive performance metric for classification tasks is the **Accuracy**.
Given a binary classification task where the neural network classifies positive and negative classes we can define the Accuracy metric as such:

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN}$$

True Positives (TP): Predicted positive and its actually positive.
True Negatives (TN): Predicted negative and its actually negative.
False Positives (FP): Predicted positive and its actually negative.
False Negative (FN): Predicted negative and its actually positive.

Accuracy is simple ratio of correctly predicted samples to the total samples and often the indication of high accuracy can be misleading especially when the class distribution of a dataset is imbalanced.
Another performance metric that is more suited for imbalanced datasets is the **F1 Score** and in order to define it we need to introduce two additional performance metrics the **Precision** and **Recall**:

- **Precision**: the ratio of correctly predicted positive observations to the total positive observations. High precision relates to the low false positive rate

$$Precision = \frac{TP}{TP + FN}$$

- **Recall**: the ratio of correctly predicted positive observations to the all observations in actual class

$$Recall = \frac{TP}{TP + FP}$$

- **F1 Score**: the weighted average of Precision and Recall. Therefore, this score takes both false positives and false negatives into account

$$F_1 \text{ Score} = \frac{2 \cdot (Recall \cdot Precision)}{Recall + Precision}$$

**Confusion matrices** are also very useful in neural network performance evaluation since they summarize all the above metrics in a visual way:



Figure 9: Confusion matrix

## 3.6   Tensors And Feature Scaling

A **tensor** is a dimensional data structure and a generalization of matrices to N-dimensional space. We can think of tensors as data containers with multiple dimensions. **Scalars**, **vectors**, and **matrices** can be represented as 0D tensors, 1D tensors and 2D tensors respectively. Data that needs a representation of higher dimensionality can use 3D, 4D .etc tensors.



Figure 10: Tensors Visual Representation

18

**Feature Scaling** is a very important procedure used to normalize the range of independent variables or features of data. In data processing, it is also known as data normalization and is generally performed during the data preprocessing step. The most common technique called min-max normalization:

$$x' = \frac{x - min(x)}{max(x) - min(x)}$$

With min-max normalization we re-scale the range of features to scale the range in [0, 1]. Here, max(x) and min(x) are the maximum and the minimum values of the feature respectively

## 3.7   Splitting to Training and Testing Dataset

Data splitting is an important aspect of neural network training procedure. The evaluation of a neural network performance must be performed on new data. In other words we use data that the network has never "seen" in order to evaluate its true performance so it is essential to avoid using the same data on training and testing. To emphasize the importance of dataset splitting we introduce the concept of **overfitting**.

**Overfitting** is when a neural network performs well on training data and poorly on new-unseen data. An overfitted neural network has learnt the patterns and the noise of the training dataset instead of the actual relationships between the features. This neural network may have a good performance on the training dataset but when a new dataset is provided as input the results are not as good , in other words the neural network does not generalize well. When we perform data splitting it is essential to not mix training samples with testing samples because the results will be misleading. Some of the techniques used for dataset splitting are:

- **Train-Test Split**:it is the simplest method of data splitting. The ratio used for the split is usually 80% training 20% testing.



Figure 11: Data split visualization

  This method has some **caveats** since some samples that share the same patterns or noise could be either on the train or the test side and this could also result in overfitting.

- **K-fold Cross Validation**:this method has many variations but the main idea is to perform multiple splits on the dataset and perform the testing on various portions of the dataset and average the results.

Figure 12: K-fold cross validation split

## 3.8 Brief Introduction to LSTMs

**LSTMs** are a special case of deep neural networks that are capable of processing sequences of data , learn long term dependencies and retain useful information from previous seen data.

An LSTM is composed by a set of **stacked units** called **cells**. An **LSTM cell** operates by using three types of gates which control the information in a sequence of data; the **forget gate , input gate** and **output gate**. We will briefly mention the role of each gate , how they communicate, operate and use the aforementioned activation functions , **sigmoid** and **tanh**.

The output of an LSTM for a specific time point is dependent on:

- The current long-term memory of the network , the **cell state**

- The output of the previous time point , the **hidden state**

- The **input data** of the current time point

In the image below we can see a high level visualization of the components of an LSTM cell.

Figure 12: Internals of a lstm cell

As we can see the main components of an LSTM cell are its three gates:

- **Forget Gate**:The forget gate is responsible for the long term memory of the lstm in the form of a **vector** and decides which piece of information is relevant to keep. In order to achieve this the previous hidden state and the new input data are fed into the forget gate and the sigmoid function outputs values close to 0 when a component of the input is deemed irrelevant and closer to 1 when relevant. The output values are then sent up and **pointwise multiplied** with the previous cell state. This pointwise multiplication means that components of the cell state which have been deemed irrelevant by the forget gate network will be multiplied by a number close to 0 and thus will have less influence on the following steps.

- **Input Gate**:The input gate takes as input the current input data and the previous hidden state and passes them through a tanh activation function and a sigmoid function in parallel. The tanh activation function outputs values close to 1 when the network estimates that this fraction of information must have impact on the outcome and close to -1 if the impact needs to be reduced.The sigmoid activation function acts as a filter that identifies which fraction of the **new memory vector** are worth retaining. This gate will output a vector of values in [0,1].Finally the outputs of the tanh and sigmoid functions are pointwise multiplied and this decides the magnitude of the new information that we use to update the long term memory of the network by performing a pointwise addition with the previous cell state

- **Output Gate**:The output gate decides the **next hidden state** or the final **prediction** of the network if the time point fed is the last of the sequence. In a similar way with the input gate the network passes the current input and the previous hidden state through a sigmoid activation function and creates a filter for the worth keeping information of the new data. Then the vector of the new cell state that contains all the information from previous timesteps and the current one is being set to a tanh activation function.The outputs of the tanh and sigmoid functions are pointwise multiplied and the result is the new hidden state.

A **sequence of stacked LSTM cells** form a LSTM neural network as illustrated bellow:



Figure 13: Stacked lstm cells forming a lstm neural net

## 3.9 Output Relations of LSTMs

Depending the task LSTMs can be configured to have various output modes regarding the number of their output. We present two of them:

- **Many-to-One**: Given a sequence of data as input, and we have to predict a single output.

Many-to-one

Figure 14: Many to one lstm

- **Many-to-Many**:Given a sequence of data as input, and we have to predict multiple outputs. In this case the size of the input can match the size of the output and this is something that we will apply in our task.



Many-to-Many

Figure 15: Many to many lstm

## 3.10   Bi-directional LSTM

Bi-directional LSTMs is a special case of an LSTM neural network. By reversing the flow of the input it can use information from the future steps of the sequence. Basically it retains an additional layer with the reversed sequence which allows the combination of information from both flows.

Figure 16: Bi-directional lstm

# 4 NeurAsp Framework Overview

In this section we will introduce the main concepts of the NeurAsp framework and in parallel we will go through the digit addition example [Yang, 2020] step by step.

## 4.1 Syntax

NeurAsp syntax is similar to ASP with a few differences that we will address in the following section.

Let's assume that a neural network $M$ takes as input a tensor $t$ and outputs a matrix in $R^{e \times n}$ where $e$ is the number of random events and $n$ the number of possible outcomes for each random event.Each row of the output matrix represents the probability distribution of each random event.The aforementioned neural network can be represented in NeurAsp as the following **neural atom**:

$$nn(m(e, t), [v_1, ..., v_n]) \ ,$$

where **nn** is a reserved keyword that denotes the neural atom , **m** is the symbolic name of the neural network $M$ and finally $[v_1, ..., v_n]$ are the **n** possible outcomes of each random event.

The task is, given a pair of digit images (MNIST) and their **sum as the label**, let a **Convolutional neural network (CNN)** learn the digit classification of the input images.

Let's define a neural atom for the digit addition task:

$$nn(digit(1, D), [1, 2, 3, 4, 5, 6, 7, 8, 9]) \ ,$$

where:

- **t=D**: the image tensor provided as input

- **e=1**: one image tensor can be classified to one digit

- **n=10**: there are 10 possible digits that the image tensor can be classified

- **possible outcomes**: [0,1,2,3,4,5,6,7,8,9]

- **m=digit**: the symbolic name of the task's neural network

- **output matrix**: a matrix in $R^{1 \times 10}$ that represents the probability distribution over the possible digits.

Since we need to apply digit addition between pairs of images we need to define the addition operation rule as such:

$$addition(D_1, D_2, N) \leftarrow digit_1(D_1), digit_2(D_2), N = D_1 + D_2 \ ,$$

In NeurAsp syntax we use the atom **img(T).** to indicate that this atom represents the input tensor T. Since we have introduced everything needed in order to formulate the digit addition task using the NeurAsp syntax we can define the NeurAsp program bellow:

$img(D_1).$
$img(D_2).$
$nn(digit(1, D), [1, 2, 3, 4, 5, 6, 7, 8, 9]) \leftarrow img(D).$
$addition(D_1, D_2, N) \leftarrow digit_1(D_1), digit_2(D_2), N = D_1 + D_2.$

## 4.2 Semantics and Internal Operations

A NeurAsp program $\Pi$ needs to be transformed to an actual ASP program called $\Pi'$. During the transformation of the program NeurAsp performs a set of actions that are needed for the framework to operate:

- **Neural atom replacement**: The neural atom is replaced by a set of rules:

  $\{m_i(t) = v_1; ...; m_i(t) = v_n\} = 1$ for $i \in \{1, ..., e\}$.

  This rule is an ASP **choice rule** in the language of **clingo** and means to choose exactly one atom between the braces. In digit addition task the choice rule would be:

  $\{digit1(0..9); digit2(0..9)\} = 1.$

  This rule would yield a set of 20 atoms (10 digits for each input image)

- **Store the atoms of choice rules in a list**: The atoms produced are stored in a python list for later usage as we will show in the training procedure.

- **Mapping of the output matrix**: Create a mapping between the output probability distribution matrix and the possible outcomes of the classification for each **input tensor**. This way each atom can be associated with a probability. In digit addition this mapping is shown in the figure below:

Atoms for all possible outcomes

| digit(0) | digit(1) | digit(2) | digit(3) | digit(4) | digit(5) | digit(6) | digit(7) | digit(8) | digit(9) |
|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|
| $p_1$ | $p_2$ | $p_3$ | $p_4$ | $p_5$ | $p_6$ | $p_7$ | $p_8$ | $p_9$ | $p_{10}$ |

Probability of each atom

Figure 16: Mapping of output probabilities

**Note** that each sample given as input to the Convolutional neural network comes as a pair of tensors which contains the numerical representation of the digit image. The neural network will output a probability distribution to its output layer for each digit image of the pair as shown in the figure below:



Figure 17: Output probabilities of a digit pair

- **Create a matrix for the computed gradient updates**: This matrix has similar mapping as the output matrix and it is used for storing the gradient updates as we will show in the training procedure.

In digit addition we provide as labels the result of the addition of two numbers and with that we can limit the choices for each digit in the given pair. For instance if we know that $D_1 + D_2 = 1$ then the possible pairs are $[D_1 = 1$ and $D_2 = 0]$ or $[D_1 = 0$ and $D_2 = 1]$.

The fact "The sum of D1 and D2 equals to 1" can be encoded as an ASP **constraint** forming an **observation**. An observation is a set of ASP constraints such as:

$$\text{:- not } addition(D_1, D_2, 1).$$

We can interpret this ASP constraint as "filter out any pair of digits that do not sum to 1" and the output would be two sets of atoms such as:

- $[digit(D_1) = 0, digit(D_2) = 1]$

- $[digit(D_1) = 1, digit(D_2) = 0]$

Those sets of atoms are called **stable models** or **interpretations I** (mentioned in Section2 **Asp and Event Calculus Background**).We can calculate the probabilities of those interpretations by using the outputs of the neural network.(As we will see in next sections where we discuss the learning procedure, the neural network outputs probabilities for each image. At first the probabilities are random but as the neural network keeps training those output probabilities are being updated).The **probability of a stable model I** is defined as the product of the probabilities of each atom they contain. For instance in our example:



Figure 18: Mapping of probabilities to atoms

$$P(I_1) = P(digit(d_1, 1)) + P(digit(d_2, 0)) = 0.3 \cdot 0.3 = 0.09$$
$$P(I_2) = P(digit(d_1, 0)) + P(digit(d_2, 1)) = 0.1 \cdot 0.1 = 0.01$$

Finally we can define the probability of an observation **O** as the sum of the stable models probabilities that satisfies **O** :

$$P(O) = \sum_{I \models O} P(I)$$

In our example that would be:

$$\text{P(:- not } addition(d_1, d_2, 1)) = \text{P(I1)} + \text{P(I2)} = 0.09 + 0.01 = 0.10$$

## 4.3  Inference with NeurAsp

Since NeurAsp has made an association between the outputs of a neural network and the atoms produced by ASP, the **reasoning about the relations among objects** recognized by the neural network is possible.

The digit addition is quite trivial and not suitable for the demonstration of the inference capabilities of the NeurAsp framework so we will use the **toycar** example provided by [Yang, 2020]. The task is that a Convolutional Neural Network can identify 5 object classes [**car,cat,person,truck,other**] in an image and we need to determine if the identified car is an **actual car** or a **toy car**. Instead of training the neural network from scratch for this specific task we can write ASP rules based on **commonsense reasoning** regarding the size relations between objects.The NeurAsp program for this task is:

nn(label(1,I,B),[car,cat,person,truck,other] :- box$(I, B, X_1, Y_1, X_2, Y_2)$.

The NeurAsp program states that the neural network takes as input the bounding box of the objects in an image and then outputs one classification label for each one of them.

In real life we know that objects can be sorted by their size and we can encode this information with the following ASP rules:

smaller(cat,person).
smaller(person,car).
smaller(person,truck).

Those rules state that the object on the left is smaller than the one in the right. We add a rule that states that if an object X is smaller than an object Y and there is an object Z that is smaller than X **then Z is also smaller than Y**.

smaller(X,Y) :- smaller(X,Z),smaller(Z,Y).

With the association of each label and its bounding box we could create a rule that by default asserts the size relationships as we defined it so far but this rule would not be sufficient to differentiate toy cars from real ones. We need to use additional information from the bounding box coordinates so we could encode the size relation of the boxes using ASP code as such:

smaller$(I, B_1, B_2)$ :- box$(I, B_1, X_1, Y_1, X_2, Y_2)$,
box$(I, B_2, X_1', Y_1', X_2', Y_2')$,
$|X_1 - X_2| \times |Y_1 - Y_2| < |X_1' - X_2'| \times |Y_1' - Y_2'|$.

Now we could make the assumption that if a car has a smaller bounding box than a person in the same picture then it is probably a toy car but this is not the case as we can see in the set of pictures below:



Figure 19: Example of inference with NeuAsp

In the left picture what we have defined so far would suffice since we have a large bounding box of a person and small ones for the cars and indeed they are toy cars , but in the picture of the right there are small bounding boxes of actual cars in the background. So there are some **exceptions regarding the size relations** that we need to encode and we will need to define a second smaller relation regarding the bounding boxes (besides the one that refers to real life sizes).

$$\text{smaller}(I, B_1, B_2) \text{ :- not } \sim\text{smaller}(I, B_1, B_2),$$
$$\text{label}(I, B_1) = L_1, \text{label}(I, B_2) = L_2, \text{smaller}(L_1, L_2).$$

The symbol $\sim$ in front of the new size relation stands for strong negation and is needed in order to assert explicit falsity and provides a way to define the exception in our size relation. The rule above can be translated as $B_1$ is smaller than $B_2$ if the exception does not hold and if $B_1$ **is smaller than $B_2$ in actual object sizes**.

In order to determine if the recognized car is a toy or not we need to apply commonsense reasoning for the bounding boxes as well. Objects in a picture appear larger if they are closer to the camera so we can use their bottom coordinates ($Y_2$) to determine which object is closer:

$$\sim\text{smaller}(I, B_2, B_1) \text{ :- box}(I, B_1, X_1, Y_1, X_2, Y_2),$$
$$\text{box}(I, B_2, X_1', Y_1', X_2', Y_2'),$$
$$Y_2 \geq Y_2',$$
$$|X_1 - X_2| \times |Y_1 - Y_2| < |X_1' - X_2'| \times |Y_1' - Y_2'|.$$

The rule above says that $B_2$ **is not smaller than** $B_1$ if $B_1$ **is closer to the camera** ($Y_2 \geq Y_2'$) and the box of $B_1$ **is smaller than** $B_2$. Then we define the rule :

$$\text{smaller}(I, B_1, B_2) \text{ :- } \sim\text{smaller}(I, B_2, B_1).$$

This rule says that $B_1$ **is smaller than** $B_2$ if $B_2$ **is not smaller than** $B_1$. Finally we can create the rule that states when a car is a toy or not:

toy$(I, B_1)$ :- label$(I, B_1) = L_1$, label$(I, B_2) = L_2$,
smaller$(I, B_1, B_2)$ ,smaller$(L_2, L_1)$.

This final rule says that the object $L_1$ with the bounding box $B_1$ is a toy if it is smaller in terms of the bounding box rules that we defined (with the exceptions of the camera distance) and also if the object $L_2$ is smaller than $L_1$ in real life.Using those rules NeurAsp is able to differentiate toy cars by applying commonsense reasoning on a Neural Network outputs.

## 4.4 Learning with NeurAsp

When we train a neural network with NeurAsp our task is to find the parameters $\theta$ that maximize the log-likelihood of the observations.

$$\theta \in argmax_\theta \sum \log(P_{\Pi(\theta)}(O))$$

The gradient of $\sum \log(P_{\Pi(\theta)}(O))$ with respect to $\theta$ is :

$$\frac{\partial \sum \log(P_{\Pi(\theta)}(O))}{\partial \theta} = \sum \frac{\partial \log(P_{\Pi(\theta)}(O))}{\partial p} \times \frac{\partial p}{\partial \theta}$$

The term $\frac{\partial p}{\partial \theta}$ can be computed through the usual neural network backpropagtion while each $\frac{\partial \log(P_{\Pi(\theta)}(O))}{\partial p}$ can be computed as shown in [Yang, 2020] , **proposition 1**:

$$\frac{\partial \log(P_{\Pi(\theta)}(O))}{\partial p} = \frac{\sum_{I:|=O} \frac{P_{\Pi(\theta)}(I))}{P_{\Pi(\theta)}(c=v))} - \sum_{I:|=O} \frac{P_{\Pi(\theta)}(I))}{P_{\Pi(\theta)}(c=v'))}}{\sum_{I:|=O} P_{\Pi(\theta)}(I))}$$

In order to explain in depth the above formula we need to break it down as shown below:

Figure 19: Formula break down

Continuing our example of digit addition we assume that at this point NeurAsp has the probability distributions for both images of the sample. In order to calculate the gradients we need the following:

- The Denominator: we need the probability of each stable model. We calculate those probabilities as the multiplication of the probabilities of each atom that each model contains as shown in the previous sub-section.

- The Numerator and its parts: NeurAsp keeps a list of all generated atoms (20 in total) as we have mentioned in sub-section **4.2 Semantics and Internal Operations** in the form of:

$$[digit(d_1,0), ..., digit(d_1,9), digit(d_2,0), ..., digit(d_2,9)]$$

  NeurAsp iterates through all the generated digit atoms and checks if they are contained in each model **(satisfied by the model)** and depending on the case it performs different updates to the gradients as we denote with the **two parts of the numerator**.
  The **first part of the numerator** indicates that the probability $p_i$ of an $atom_i$ tends to increase for every stable model that contains this particular atom.
  On the other hand the **second part of the numerator** indicates that the probability $p_i$ of an $atom_i$ tends to decrease for each stable model that does not contain this particular atom.
  The **value** of each increase or decrease its dependent on the probability of each atom and the probability of each stable model of the observation as they are defined in the previous sub-section.

  In order to elaborate more on the calculation of **gradient updates** we provide a step by step example for some of the generated output atoms:

Assuming that we have the following output probabilities and stable models of the observation:

**Probabilities of digit1**

| 0.1 | 0.3 | 0.1 | 0.1 | 0.05 | 0.03 | 0.02 | 0.1 | 0.1 | 0.1 |
|-----|-----|-----|-----|------|------|------|-----|-----|-----|

**Probabilities of digit2**

| 0.3 | 0.1 | 0.1 | 0.05 | 0.03 | 0.02 | 0.1 | 0.1 | 0.1 | 0.1 |
|-----|-----|-----|------|------|------|-----|-----|-----|-----|

**Stable Models**

| Model1 (I1) | digit(d1,1) | digit(d2,0) |
|-------------|-------------|-------------|
| Model2 (I2) | digit(d1,0) | digit(d2,1) |

and the computed probabilities of the stable models:

$$P(I_1) = P(digit(d_1, 1)) + P(digit(d_2, 0)) = 0.3 \cdot 0.3 = 0.09$$
$$P(I_2) = P(digit(d_1, 0)) + P(digit(d_2, 1)) = 0.1 \cdot 0.1 = 0.01$$

NeurAsp will need to **store the computed gradients** before proceeding with the backpropagation so for this reason there is a matrix initialized with zeroes. This matrix has a mapping with all possible outcomes ( in our case 9 outcomes) and will be populated with the gradient updates as we will describe.

**Gradient Update Matrix**

| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
|---|---|---|---|---|---|---|---|---|---|

NeurAsp will iterate though all 20 generated atoms and for each of the two stable models starting with $digit(d_1, 0)$ will do the following:

- $digit(d_1, 0)$ is not in model1:$[digit(d_1, 1), digit(d_2, 0)]$

  In this model the atom is not contained since for $d_1$ the atom $digit(d_1, 1)$ is present with probability $p_1 = 0.3$ so we will have a **negative gradient update** of value $\frac{P(model_1)}{p_1} = \frac{0.09}{0.3}$ to be added in the 0 index of gradient update matrix.

- $digit(d_1, 0)$ is in model2:$[digit(d_1, 0), digit(d_2, 1)]$

In this model the atom with probability $p_0 = 0.1$ is contained so we will have a **positive gradient update** of value $\frac{P(model_1)}{p_0} = \frac{0.09}{0.1}$ to be added in the 0 index of gradient update matrix.

Continuing with $digit(d_1, 1)$:

- $digit(d_1, 1)$ is in model1:$[digit(d_1, 1), digit(d_2, 0)]$

  In this model the atom is contained with probability $p_1 = 0.3$ so we will have a **positive gradient update** of value $\frac{P(model_1)}{p_1} = \frac{0.09}{0.3}$ to be added in the 1st index of gradient update matrix .

- $digit(d_1, 1)$ is not in model1:$[digit(d_1, 0), digit(d_2, 1)]$

  In this model the atom is not contained since for $d_1$ the atom $digit(d_1, 0)$ is present with probability $p_0 = 0.1$ so we will have a **negative gradient update** of value $\frac{P(model_1)}{p_0} = \frac{0.09}{0.1}$ to be added in the 1st index of gradient update matrix.

We will provide a last iteration with $digit(d_1, 2)$:

- $digit(d_1, 2)$ is not in model1:$[digit(d_1, 1), digit(d_2, 0)]$

  In this model the atom is not contained since for $d_1$ the atom $digit(d_1, 1)$ is present with probability $p_1 = 0.3$ so we will have a **negative gradient update** of value $\frac{P(model_1)}{p_1} = \frac{0.09}{0.3}$ to be added in the 2nd index of gradient update matrix.

- $digit(d_1, 2)$ is not in model1:$[digit(d_1, 0), digit(d_2, 1)]$

  In this model the atom is not contained since for $d_1$ the atom $digit(d_1, 0)$ is present with probability $p_0 = 0.1$ so we will have a **negative gradient update** of value $\frac{P(model_1)}{p_0} = \frac{0.09}{0.1}$ to be added in the 2nd index of gradient update matrix.

Any further iteration for digit $d_1$ will result to a negative update. Again when the loop reaches atoms regarding digit $d_2$ it will do the same respective updates. When all atoms are checked the gradient update matrix is populated and backpropagation is initiated.

# 5 Human Activity Recognition with NeurAsp

## 5.1 Task Description

NeurAsp showed promising results compared to traditional deep neural network methods regarding the digit addition task. The purpose of this thesis is to investigate the capabilities of the NeurAsp framework with a more challenging task such as the recognition of the interactions between two persons in a video.

Given a set of video frames where persons are performing a set of actions (**simple events**) such as walking , running , remain idle etc we will train a neural network that will recognize the interaction between pairs of persons and categorize them to **complex events** classes such as meeting , moving together etc. Using as a "teacher" a logical theory that provides definitions of the complex events in terms of the simple ones. That is, the task is to train the Neural Network to predict simple events in a way that optimizes the predictive performance of the complex event rules when they consume such simple event predictions. The simple events here will be treated as latent concepts, similarly to single-digit labels in the MNIST addition task in previous section

For this task we will need to create suitable input tensors that contain a sequence of features for a given number of frames (timepoints). Since sequential data is involved we have chosen to train a **Bi-directional LSTM** that will output classification results for each timepoint given as input.

## 5.2 CAVIAR Dataset

The dataset used for this thesis is the CAVIAR dataset. We are using the first section of video clips recorded in the entrance of the INRIA Labs at Grenoble, France[3]. All clips are recorded at 25 frames per second and include actions performed by individual persons or by groups of persons. Each frame is annotated in **XML format** with a label and bounding boxes that bound the people involved.We refer to actions that involve individual people as **simple events** and those events are labeled by the annotators using the following labels:

- **inactive**

- **active**

- **walking**

- **running**

Actions that groups of persons perform are referred as **complex events** and the following labels are used by the annotators:

- **joining**

- **meeting**

- **moving**

- **split up**

- **leaving object**

- **fighting**

The **features** provided for each frame are:

---

[3]https://homepages.inf.ed.ac.uk/rbf/CAVIARDATA1

- **orientation**: the orientation of the observed person

- **cx**: x coordinate of the bounding box center

- **cy**: y coordinate of the bounding box center

- **w**: width of the bounding box

- **h** : height of the bounding box

## 5.3    Data Pre-processing

In order to perform our experiments, the CAVIAR dataset needed to be processed. Each video is described by an XML file that contains information about what is happening in each frame regarding each person or group of persons. We need to create suitable input and output tensors for our **LSTM** and this is achieved by performing a set of data pre-processing steps as described below.

### 5.3.1    XML parsing

A python script extracts for each video all the necessary information from the XML file and this information is stored to a python dictionary. This dictionary contains information about all the videos provided by CAVIAR, organized by each video source. For each frame in every video we keep the following information:

- all **visible persons** ids in each frame (each person has its own id)

- features of each visible person (orientation, center coordinates of bounding boxes and their corresponding width and height) and its simple event label

- person ids that are involved in a complex event and the complex event label

- all **possible pairs** of person ids that are visible in the frame regardless of whether they are performing a complex event or not. If the pair is not performing a complex event as dictated by the complex event annotation their complex event label is annotated as "no interaction"

An example for one frame with 4 visible persons is provided below where the text colors match the bounding box colors in Figure 20:

Figure 20: Sample image of individuals and their bounding boxes

The information from the annotators for this frame is:

**Simple Events**:

- Person1: walking

- Person2: walking

- Person3: inactive

- Person4: walking

**Complex Events** annotated as a group action in XML file:

Person1 and Person2 : moving

Given this information we can produce the following **pairs** and also annotate the complex event of each **non participating pair** as a "no interaction" complex event:

- Person1 - Person2: moving

- Person1 - Person3: no interaction

- Person1 - Person4: no interaction

- Person2 - Person3: no interaction

- Person2 - Person4: no interaction

- Person3 - Person4: no interaction

A person's id is an integer starting from 1. In our example we have 4 persons with ids from 1 to 4 and also in another video we may have the same ids as numbers but they refer to other persons so we need to find a way to differentiate them as pairs.

The solution comes if we create **pair ids** that contain also the source file name. For instance if the source file is called wk1gt the pair ids in our examples would be : wk1gt_1-2 , wk1gt_1-3 , wk1gt_1-4 etc.

### 5.3.2 Sliding window on video frames

A time series classification task requires specific pre-process in order to associate signal data (in our case video frames) with the activity labels. A straightforward data partitioning approach is to divide the input data into time windows where a given window may have a few timepoints of observation data. This data preparation is called **sliding window** and it has many variations [Laguna, 2011]. After we have stored our data by **distinct pairs** we can create data chunks on a given time interval **T** of **consecutive frames**. Since in the CAVIAR dataset the number of samples of consecutive frames is somehow limited we apply the **sliding window technique** with **coverage C** in order to increase the number of our samples. Let's assume that we have the following set of 18 consecutive frames for one person of the pair and we need to create samples with 6 timepoints (window size T=6):

| Frame Num | FEATURES | | | | |
|---|---|---|---|---|---|
| | orientation | cx | cy | h | w |
| 0 | 131 | 184 | 47 | 32 | 34 |
| 1 | 131 | 184 | 47 | 32 | 34 |
| 2 | 131 | 184 | 47 | 32 | 34 |
| 3 | 131 | 184 | 47 | 32 | 34 |
| 4 | 131 | 184 | 47 | 32 | 34 |
| 5 | 131 | 184 | 47 | 32 | 34 |
| 6 | 131 | 184 | 47 | 32 | 34 |
| 7 | 131 | 184 | 47 | 32 | 34 |
| 8 | 131 | 184 | 46 | 31 | 35 |
| 9 | 131 | 184 | 46 | 31 | 35 |
| 10 | 131 | 185 | 46 | 32 | 35 |
| 11 | 131 | 185 | 46 | 32 | 35 |
| 12 | 131 | 185 | 46 | 33 | 35 |
| 13 | 131 | 186 | 46 | 32 | 33 |
| 14 | 131 | 186 | 46 | 32 | 33 |
| 15 | 131 | 186 | 46 | 32 | 33 |
| 16 | 131 | 187 | 45 | 31 | 35 |
| 17 | 132 | 187 | 47 | 33 | 34 |

Figure 21: Sample of features for each time point

If we create samples of size T=6 with a simple segmentation of the frames we will end up with 18 / 6 = 3 samples (Figure 4)

| Frame Num | FEATURES | | | | |
|---|---|---|---|---|---|
| | orientation | cx | cy | h | w |
| 0 | 131 | 184 | 47 | 32 | 34 |
| 1 | 131 | 184 | 47 | 32 | 34 |
| 2 | 131 | 184 | 47 | 32 | 34 |
| 3 | 131 | 184 | 47 | 32 | 34 |
| 4 | 131 | 184 | 47 | 32 | 34 |
| 5 | 131 | 184 | 47 | 32 | 34 |
| 6 | 131 | 184 | 47 | 32 | 34 |
| 7 | 131 | 184 | 47 | 32 | 34 |
| 8 | 131 | 184 | 46 | 31 | 35 |
| 9 | 131 | 184 | 46 | 31 | 35 |
| 10 | 131 | 185 | 46 | 32 | 35 |
| 11 | 131 | 185 | 46 | 32 | 35 |
| 12 | 131 | 185 | 46 | 33 | 35 |
| 13 | 131 | 186 | 46 | 32 | 33 |
| 14 | 131 | 186 | 46 | 32 | 33 |
| 15 | 131 | 186 | 46 | 32 | 33 |
| 16 | 131 | 187 | 45 | 31 | 35 |
| 17 | 132 | 187 | 47 | 33 | 34 |

Figure 22: Partitioning without sliding windows

By applying the **sliding window technique** with window size T=6 and coverage C=50% we will end up with 5 samples (Figure 5)

### 5.3.3   Ground Truth Generation on Complex Events with ASP

In this section we will present the rules used to detect the complex events of the task as proposed in WOLED [Katzouris, 2020]. The CAVIAR dataset provides ground truth for both simple and complex events but in order to train the neural network to recognise simple events in a way that maximizes the predictive performance of the rules we need to "extract" the ground truth for complex events by following the described procedure bellow.

The **Event Calculus** axioms are **domain independent** axioms and describe the **law of inertia** where an event holds at time point $T + 1$ if it is initiated at the previous time point $T$, or if it holds at $T$ and it is not terminated at $T$.

In ASP we can define the law of inertia with these rules:

1. holdsAt(F,Ts+1) :- initiatedAt(F,Ts), fluent(F), next(Ts,Te).
2. holdsAt(F,Te) :- holdsAt(F,Ts), not terminatedAt(F,Ts), fluent(F), next(Ts,Te).

| Frame Num | FEATURES | | | | |
|---|---|---|---|---|---|
| | orientation | cx | cy | h | w |
| 0 | 131 | 184 | 47 | 32 | 34 |
| 1 | 131 | 184 | 47 | 32 | 34 |
| 2 | 131 | 184 | 47 | 32 | 34 |
| 3 | 131 | 184 | 47 | 32 | 34 |
| 4 | 131 | 184 | 47 | 32 | 34 |
| 5 | 131 | 184 | 47 | 32 | 34 |
| 6 | 131 | 184 | 47 | 32 | 34 |
| 7 | 131 | 184 | 47 | 32 | 34 |
| 8 | 131 | 184 | 46 | 31 | 35 |
| 9 | 131 | 184 | 46 | 31 | 35 |
| 10 | 131 | 185 | 46 | 32 | 35 |
| 11 | 131 | 185 | 46 | 32 | 35 |
| 12 | 131 | 185 | 46 | 33 | 35 |
| 13 | 131 | 186 | 46 | 32 | 33 |
| 14 | 131 | 186 | 46 | 32 | 33 |
| 15 | 131 | 186 | 46 | 32 | 33 |
| 16 | 131 | 187 | 45 | 31 | 35 |
| 17 | 132 | 187 | 47 | 33 | 34 |

Sliding window1, Sliding window2, Sliding window3, Sliding window4, Sliding window5

Figure 23: Partitioning with sliding windows

where:

fluent/1 predicate: captures properties that persist in time

next/2 predicate: defined as $\mathrm{next}(T_1, T_2) \coloncolon T_2 = T_1 + 1$

In this thesis we focus on two types of complex events: the **moving** and **meeting** complex event. **Table** 2 and **Table** 3 contain the rules that define the **moving complex event** and **meeting complex event** respectively.

The moving Complex Event consists of the following rules:

1. initiatedAt(moving(X0,X1),T) :- happensAt(walking(X0),T),
   happensAt(walking(X1),T),orientationMove(X0,X1,T),
   close(X0,X1,D,T).
2. terminatedAt(moving(X0,X1),T) :- happensAt(walking(X0),T),
   far(X0,X1,D,T).
3. terminatedAt(moving(X0,X1),T) :- happensAt(walking(X1),T),
   far(X0,X1,D,T).
4. terminatedAt(moving(X0,X1),T) :- happensAt(active(X0),T),
   happensAt(active(X1),T).
5. terminatedAt(moving(X0,X1),T) :- happensAt(active(X0),T),
   happensAt(inactive(X1),T).
6. terminatedAt(moving(X0,X1),T) :- happensAt(active(X1),T),
   happensAt(inactive(X0),T).
7. terminatedAt(moving(X0,X1),T) :- happensAt(running(X0),T),
   person(X1).
8. terminatedAt(moving(X0,X1),T) :- happensAt(running(X1),T),
   person(X0).
9. terminatedAt(moving(X0,X1),T) :- happensAt(disappear(X0),T),
   person(X1).
10. terminatedAt(moving(X0,X1),T) :- happensAt(disappear(X1),T),
    person(X0).

Table 2: Definition of moving complex event

The first rule dictates the moving event between two persons X0, X1 is initiated at time T if both X0 and X1 are walking at T and their euclidean distance at T (close/4 predicate) is less than a threshold D (which is explicitly provided at runtime as a set of choices/bins for D) and their difference in orientation is such that the two persons are moving almost in parallel at T (orientationMove/3 predicate). The rest of the rules define the termination of this event. For instance (rules no2 and no3) dictate that the moving event is terminated at time T if the eulcidean distance at T (far/4 predicate) is greater that the thershold D. The moving event also terminates at time T if the two persons X0,X1 cease to perform the walking simple event concurrently at T (X0 is walking and X1 is inactive) as dictated from rules 3-6. Finally the moving event terminates if one of the persons X0, X1 starts running (rules no7 and no9) or one of them is leaving the monitored area (disappear/2 predicate) at time T.

The meeting Complex Event consists of the following rules:

1. initiatedAt(meeting(X0,X1),T):-happensAt(active(X0),T),
   happensAt(active(X1),T),close(X0,X1,D,T).
2. initiatedAt(meeting(X0,X1),T) :- happensAt(active(X0),T),
   happensAt(inactive(X1),T),close(X0,X1,D,T).
3. initiatedAt(meeting(X0,X1),T) :- happensAt(inactive(X0),T),
   happensAt(active(X1),T),close(X0,X1,D,T).
4. initiatedAt(meeting(X0,X1),T) :- happensAt(inactive(X0),T),
   happensAt(inactive(X1),X2),close(X0,X1,D,X2).
5. terminatedAt(meeting(X0,X1),T) :- happensAt(running(X0),T),
   person(X1).
6. terminatedAt(meeting(X0,X1),T) :- happensAt(running(X1),T),
   person(X0).
7. terminatedAt(meeting(X0,X1),T) :- happensAt(disappear(X0),T),
   person(X1).
8. terminatedAt(meeting(X0,X1),T) :- happensAt(disappear(X1),T),
   person(X0).
9. terminatedAt(meeting(X0,X1),T) :- happensAt(walking(X0),T),
   far(X0,X1,D,T).
10. terminatedAt(meeting(X0,X1),T) :- happensAt(walking(X1),T),
    far(X0,X1,D,T).

Table 3: Definition of meeting complex event

The first 4 rule dictate the meeting event between two persons X0, X1 is initiated at time T if both X0 and X1 are active or inactive at T and their euclidean distance at T (close/4 predicate) is less than a threshold D. The rest of the rules define the termination of this event. The meeting event terminates if one of the persons X0, X1 starts running or walking or one of them is leaving the monitored area (disappear/2 predicate) at time T.

Given the provided definitions for the two complex events we can create helper rules that define "types" of domain variables such as the complex events and express them by using the domain independent axioms of Event Calculus in the form of fluents that are subjects to inertia.

---

1. fluent(moving(X,Y)) :- person(X), person(Y), X != Y.
2. fluent(meeting(X,Y)) :- person(X), person(Y), X != Y.
3. time(T) :- happensAt(disappear(_),T).
4. time(T) :- happensAt(appear(_),T).
5. time(T) :- happensAt(active(_),T).
6. time(T) :- happensAt(inactive(_),T).
7. time(T) :- happensAt(walking(_),T).
8. time(T) :- happensAt(running(_),T).
9. time(T) :- coords(_,_,_,T).
10. time(T) :- orientation(_,_,T).
11. person(X) :- happensAt(disappear(X),_).
12. person(X) :- happensAt(appear(X),_).
13. person(X) :- happensAt(active(X),_).
14. person(X) :- happensAt(inactive(X),_).
15. person(X) :- happensAt(walking(X),_).
16. person(X) :- happensAt(running(X),_).

---

Table 4: Helper rules that define "types" of domain variables

The first two rules define the complex event as fluents. (the fluent atom is used in the inertia law ASP definition).
The happensAt atom is used to associate:

- time : when the event started or ended

- event: what event started or ended (events such as walking , inactive , active and running)

- person: the id of the person that performs the event that started or ended

Given a **pair of persons** in a **time range** with the rules mentioned above we can determine if those persons are moving together , meeting or having no interaction at all. (there is no rule for no interaction event because if those two persons are not *moving together* or *meeting* the ASP program will not output anything about the given time range so it is assumed that there is no interaction at all).
Since we have parsed and organized our data by unique pairs we can create separate asp programs for each individual pair and retrieve the time intervals in which the complex events hold continuously. Each ASP program contains:

- The **rule of inertia** and the **complex event definitions** as **background knowledge**.

- The coordinates , orientation , the ids of the persons and the frame number encoded as asp atoms. (coords , orientation etc)

- The **ground truth** of the simple events performed by each person in the pair encoded as **happensAt** atoms.

As mentioned our goal is to maximize the predictive performance of the defined complex event rules (**Table** 2 and **Table** 3) through neuro-symbolic training. We treat those rules as the "golden standard" and by providing as input the ground truth on simple events as given from CAVIAR we extract the ground truth for complex events.
After running each ASP program we obtain time intervals where the **complex event holds**. In our experiments we will use those intervals as **ground truth for the complex events**.

### 5.3.4 Stable Model generation for time windows

In digit addition [Section 4.2] we provided **observations** regarding the result of the addition between two digits and with ASP we generated the possible stable models. By applying the same logic we provide observations regarding the complex event that two persons are performing and we will generate stable models on the **simple events** that those two persons may perform **individually**.

Let's assume that we have two persons that are performing a **meeting** complex event in the time range **73-77**. For this time range of **5 frames** we create the following choice rules that will generate all the **possible combinations** of the simple event that a person may perform:

1{simpleEvent(T,person_1,active); simpleEvent(T,person_1,inactive); simpleEvent(T,person_1,walking); simpleEvent(T,person_1,running)}1 :- T=72..78.

1{simpleEvent(T,person_2,active); simpleEvent(T,person_2,inactive); simpleEvent(T,person_2,walking); simpleEvent(T,person_2,running)}1 :- T=72..78.

**Note** that in the time range of the choice rules we added one frame to the start and one to the end ( the original time range was 73-77 but in the choice rules is 72-78) where we assume that the complex event does not hold. This addition is needed to trigger the rules that use the law of inertia.

Given that we know that those two persons are meeting in the range of 73-77 we can create the following set of **observations** with the addition of set of rules where the complex event does not hold (set1 / set7).

**Note** that these observations are encoded via ASP constraints that enforce their satisfaction by an ASP program used for reasoning:

1. :- holdsAt(meeting(person_1,person_2),72).
   :- holdsAt(meeting(person_2,person_1),72).
2. :- not holdsAt(meeting(person_1,person_2),73).
   :- not holdsAt(meeting(person_2,person_1),73).
3. :- not holdsAt(meeting(person_1,person_2),74).
   :- not holdsAt(meeting(person_2,person_1),74).
4. :- not holdsAt(meeting(person_1,person_2),75).
   :- not holdsAt(meeting(person_2,person_1),75).
5. :- not holdsAt(meeting(person_1,person_2),76).
   :- not holdsAt(meeting(person_2,person_1),76).
6. :- not holdsAt(meeting(person_1,person_2),77).
   :- not holdsAt(meeting(person_2,person_1),77).
7. :- holdsAt(meeting(person_1,person_2),77).
   :- holdsAt(meeting(person_2,person_1),78).

These observation rules contain in their body the holdsAt atom that is used to define the complex event definitions and will **filter out** any generated model that violates any of those rules leaving behind only the models that comply. Essentially, these constraints enforce that the complex events ground truth that corresponds to the rules of **Table** 2 and **Table** 3 is satified.
In order to generate models that contain simpleEvent atoms we add the following rules that map the happensAt atoms to simpleEvent atoms:

1 .happensAt(active(X),T) :- simpleEvent(T,X,active).
2 .happensAt(inactive(X),T) :- simpleEvent(T,X,inactive).
3 .happensAt(walking(X),T) :- simpleEvent(T,X,walking).
4. happensAt(running(X),T) :- simpleEvent(T,X,running).

Finally we provide atoms such as **coordinates , orientation** and **visibility** extracted from the XML files and we create separate ASP programs for each time window. **Here we will emphasize the main difference with the previous phase of ASP code generation**. In this phase we do not provide any information about the simple events performed since we need to do the **"reverse procedure"**. Specifically:

- In the first phase we generated models that contained atoms with time intervals (frames) where a complex event holds **given the simple events provided by the XML files**. We did this in order to acquire the deductive consequences of our "teacher" rules (the meeting/moving definitions) as the complex event ground truth that will be used in the neuro-symbolic training.

- In this case **given that the complex event holds** we generate models that contain atoms with time intervals where the two involved persons perform a **simple event individually**. Therefore, in this phase we are abductively explaining the complex event ground truth (encoded as a as set of constraints) in terms of the simple event labels. Note that the latter are the "latent concept" labels that the Neural Network needs to learn how to predict using the complex event rules as a "teacher". We are thus following the standard NeurASP approach here, properly adapted to the particular use case.

In our example ASP will generate the following models:

| Model 1 | Model 2 | Model 3 | |
|---|---|---|---|
| simpleEvent(72,person_1,running) simpleEvent(72,person_2,walking) | simpleEvent(72,person_1,running) simpleEvent(72,person_2,active) | simpleEvent(72,person_1,running) simpleEvent(72,person_2,walking) | start frame 72 |
| simpleEvent(73,person_1,active) simpleEvent(73,person_2,active) | simpleEvent(73,person_1,inactive) simpleEvent(73,person_2,active) | simpleEvent(73,person_1,inactive) simpleEvent(73,person_2,active) | |
| simpleEvent(74,person_1,active) simpleEvent(74,person_2,active) | simpleEvent(74,person_1,inactive) simpleEvent(74,person_2,active) | simpleEvent(74,person_1,inactive) simpleEvent(74,person_2,active) | |
| simpleEvent(75,person_1,active) simpleEvent(75,person_2,active) | simpleEvent(75,person_1,inactive) simpleEvent(75,person_2,active) | simpleEvent(75,person_1,inactive) simpleEvent(75,person_2,active) | Intermediate frames |
| simpleEvent(76,person_1,active) simpleEvent(76,person_2,active) | simpleEvent(76,person_1,inactive) simpleEvent(76,person_2,active) | simpleEvent(76,person_1,inactive) simpleEvent(76,person_2,active) | |
| simpleEvent(77,person_1,active) simpleEvent(77,person_2,active) | simpleEvent(77,person_1,inactive) simpleEvent(77,person_2,active) | simpleEvent(77,person_1,inactive) simpleEvent(77,person_2,active) | |
| simpleEvent(78,person_1,running) simpleEvent(78,person_2,running) | simpleEvent(78,person_1,walking) simpleEvent(78,person_2,running) | simpleEvent(78,person_1,running) simpleEvent(78,person_2,running) | end frame 78 |

Figure 24: stable models

ASP in fact will generate a lot more stable models. The fact that the initiation and the termination time-point of a complex event requires to not hold, ASP will assign all possible combinations of simple events to those time-points. In order to tackle this we made the assumption that once a complex event is not re-initiated within an interval, the generated simple events within this interval match those of the initiating time-point. This simplifying assumption is justified by the fact that it is often natural to expect that the conditions that initiate a durative event are likely to re-occure for as long as the event holds. Given the assumption we made, we are interested only in the **intermediate frames** that the complex event holds so we discard the start frame and the end frame. Note that we need to keep the **unique stable models** of the **intermediate frames** and in our example **model2** and **model3** are the same. So in the time range

of our interest we keep only two instances of stable models.



Figure 25: Unique instances of stable models

### 5.3.5   Input and Label Tensors

In this thesis we are evaluating three different experimental setups. All experimental setups are presented and described in following sections but we will make a short introduction here since each setup requires different input tensors and data structures:

- **Only Deep Learning methods (setup1)**: this setup involves only traditional deep learning methods without the intervention of logic.

- **Training on simple events / Infer with NeurAsp (setup2)**: this method lets the neural network train with the usual back-propagation using the ground truth as provided by the CAVIAR dataset.

- **NeurAsp training and inference (setup3)**: this method uses the NeurAsp framework for training and inference.

Before we begin the description of the input and label tensors we remind the **5 features** provided by the CAVIAR Dataset: [**orientation, xc ,xy ,h ,w**].

Since neural network operations are numerical, tensors must only contain numbers. In order to compare the predictions made by a neural network with non numeric target values such as in our task we need to create tensors with numbers that map with each **target label**. So we have chosen the following **label**

**encoding** for complex and simple events.

| SIMPLE EVENT MAPPING | |
|---|---|
| EVENT | NUMBER |
| active | 0 |
| inactive | 1 |
| running | 2 |
| walking | 3 |

| COMPLEX EVENT MAPPING | |
|---|---|
| EVENT | NUMBER |
| no interaction | 0 |
| meeting | 1 |
| moving | 2 |

Figure 26: Label encoding for simple and complex events

After all the steps of data pre-processing that we described are completed ,the following **three sets** of **data structures and tensors** (one set for each method) are produced and we provide a visual representation of each one for window size T=6:

- **Set1 of tensors (Only Deep Learning methods)**

  - **Input tensors of size Tx11**:this tensor contains the features of two persons concatenated row wise for each time point. We added the euclidean distance between the coordinates of the two persons involved as an additional feature. ( 5 features for each person plus the distance feature).

  - **Label tensors with complex events of size Tx1**: this tensor contains the encoded labels for the complex events for each time step as generated from the ASP programs.

| Frame Num | PERSON1 FEATURES | | | | | PERSON2 FEATURES | | | | | Frame Num | COMPLEX EVENT |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | orientation | cx | cy | h | w | orientation | cx | cy | h | w | | LABELS |
| 0 | 131 | 184 | 47 | 32 | 34 | 135 | 182 | 45 | 32 | 34 | 0 | meeting |
| 1 | 131 | 184 | 47 | 32 | 34 | 135 | 182 | 45 | 32 | 34 | 1 | meeting |
| 2 | 131 | 184 | 47 | 32 | 34 | 135 | 182 | 45 | 32 | 34 | 2 | meeting |
| 3 | 131 | 184 | 47 | 32 | 34 | 135 | 182 | 45 | 32 | 34 | 3 | meeting |
| 4 | 131 | 190 | 53 | 35 | 39 | 135 | 185 | 49 | 34 | 38 | 4 | moving |
| 5 | 131 | 182 | 56 | 38 | 39 | 135 | 189 | 52 | 35 | 39 | 5 | moving |

Row-wise feature concatenation    Euclidean distance    Label encoding

**INPUT TENSOR 6x11**

| 131 | 184 | 47 | 32 | 34 | 135 | 182 | 45 | 32 | 34 | 2.82 |
|---|---|---|---|---|---|---|---|---|---|---|
| 131 | 184 | 47 | 32 | 34 | 135 | 182 | 45 | 32 | 34 | 2.82 |
| 131 | 184 | 47 | 32 | 34 | 135 | 182 | 45 | 32 | 34 | 2.82 |
| 131 | 184 | 47 | 32 | 34 | 135 | 185 | 49 | 34 | 38 | 2.82 |
| 131 | 190 | 53 | 35 | 39 | 135 | 185 | 49 | 34 | 38 | 6.4 |
| 131 | 182 | 56 | 38 | 39 | 135 | 189 | 52 | 35 | 39 | 8.06 |

**LABEL TENSOR 6x1**

| 1 |
|---|
| 1 |
| 1 |
| 1 |
| 2 |
| 2 |

Figure 27: Tensors used in setup1

- **Set2 of tensors (Training on simple events / Infer with NeurAsp)**

  - **Input tensors of size Tx5**: Each one of these tensors contain features of one person.
  - **Label tensors with simple events of size Tx1**: this tensor contains the labels on simple events **as given from the XML files**.
  - **Lists of atoms**: For each sample tensor there is a list that contains coordinates, orientation and visibility atoms that match their respective frames.
  - **Label tensors with complex events of size Tx1**: this tensor contains the labels on complex events for each time step as generated from the ASP programs.

| PERSON1 6x5 Tensor | | | | |
|---|---|---|---|---|
| 131 | 184 | 47 | 32 | 34 |
| 131 | 184 | 47 | 32 | 34 |
| 131 | 184 | 47 | 32 | 34 |
| 131 | 184 | 47 | 32 | 34 |
| 131 | 190 | 54 | 35 | 39 |
| 131 | 182 | 56 | 38 | 39 |

| PERSON2 6x5 Tensor | | | | |
|---|---|---|---|---|
| 135 | 182 | 45 | 32 | 34 |
| 135 | 182 | 45 | 32 | 34 |
| 135 | 182 | 46 | 33 | 34 |
| 135 | 182 | 45 | 32 | 34 |
| 135 | 185 | 49 | 34 | 38 |
| 135 | 189 | 52 | 35 | 39 |

**List of ASP atoms extracted from XML files**

orientation(ID,FrameNum)
coords(ID,xc,xy,FrameNum)
visible(ID,FrameNum)

Ground Truth on simple events as given from the xmls

| PERSON1 6x1 Label Tensor | |
|---|---|
| Simple Event | Encoding |
| Inactive | 1 |
| Inactive | 1 |
| Inactive | 1 |
| Active | 0 |
| Walking | 3 |
| Walking | 3 |

| PERSON2 6x1 Label Tensor | |
|---|---|
| Simple Event | Encoding |
| Inactive | 1 |
| Active | 0 |
| Inactive | 1 |
| Active | 0 |
| Walking | 3 |
| Walking | 3 |

| LABEL TENSOR 6x1 | |
|---|---|
| Simple Event | Encoding |
| Meeting | 1 |
| Meeting | 1 |
| Meeting | 1 |
| Meeting | 1 |
| Moving | 2 |
| Moving | 2 |

Figure 28: Tensors used in setup2

- **Set3 (Training and Infer with NeurAsp)**

  - **Tuple of input tensors of size Tx5**: The tensors of the two persons involved in the complex event are stored in a tuple of size 2.

  - **Lists of stable models**: Each tuple is accompanied with a list that contains the stable models produced from ASP.

  - **Lists of atoms**: For each tuple of tensors there is a list that contains coordinates , orientation and visibility atoms that match their respective frames.

  - **Label tensors with complex events of size Tx1**: this tensor contains the labels on complex events for each time step as generated from the ASP programs.

**PERSON1 6x5 Tensor**
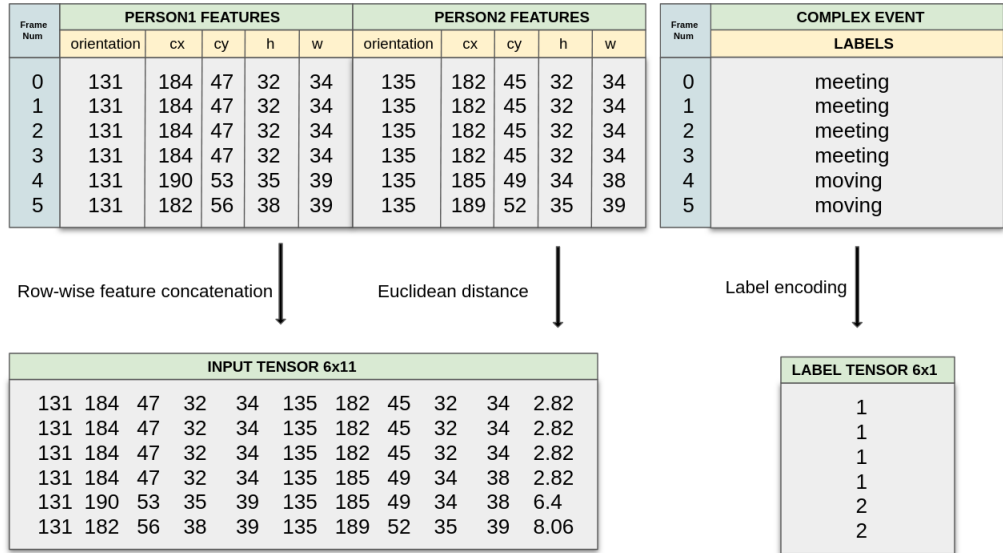
| 131 | 184 | 47 | 32 | 34 |
| 131 | 184 | 47 | 32 | 34 |
| 131 | 184 | 47 | 32 | 34 |
| 131 | 184 | 47 | 32 | 34 |
| 131 | 190 | 54 | 35 | 39 |
| 131 | 182 | 56 | 38 | 39 |

**PERSON2 6x5 Tensor**

| 135 | 182 | 45 | 32 | 34 |
| 135 | 182 | 45 | 32 | 34 |
| 135 | 182 | 46 | 33 | 34 |
| 135 | 182 | 45 | 32 | 34 |
| 135 | 185 | 49 | 34 | 38 |
| 135 | 189 | 52 | 35 | 39 |

**List of ASP atoms extracted from XML files**

orientation(ID,FrameNum)
coords(ID,xc,xy,FrameNum)
visible(ID,FrameNum)

**List of Unique stable models**

Model1: simpleEvent(person1,active,0),simpleEvent(person2,active,0),...
Model2: simpleEvent(person1,inactive,0),simpleEvent(person2,active,0),...
Model3: simpleEvent(person1,active,0),simpleEvent(person2,inactive,0),...
....

ModelN: simpleEvent(person1,inactive,0),simpleEvent(person2,active,0),...

**LABEL TENSOR 6x1**

| Simple Event | Encoding |
| --- | --- |
| Meeting | 1 |
| Meeting | 1 |
| Meeting | 1 |
| Meeting | 1 |
| Moving | 2 |
| Moving | 2 |

Figure 29: Tensors used in setup3

### 5.3.6 CAVIAR Dataset Cross Validation Split

After we have obtained the desired format of our tensors we selected **38 unique pairs of persons** for our Dataset that perform complex events. Some pairs were omitted since they performed only the "no interaction" complex event in the entire frame sequence and this class has already the majority of the samples as we can see in the following distribution diagram:

**All Complex Events Sample Distribution**

- no_interaction: 9864
- moving: 4362
- meeting: 2814

Figure 30: Distribution of complex events

As we can observe, our dataset is heavily **imbalanced**. In order to have valid and robust results we need to perform a form of a **stratified k-fold cross validation**. The term **stratified** implies that we must retain the distribution ratio in our folds. We consider each pair to be **independent** from the other

so the folds that are used for testing contain pairs that are not present in the training dataset each time.

**Note** that the number of complex events performed by each unique pair are not evenly distributed so in order to keep the distribution ratio both in the training and in the testing segments, the number of unique pairs in each fold varies but the distribution ratio is the same.

For example let's assume that we have 10 unique pairs with the following distribution:

**All samples**

- no interaction: 180 samples

- meeting: 50 samples

- moving: 80 samples

We split to 80% training and 20% testing:

**Training**

- no interaction: 144 samples

- meeting: 40 samples

- moving: 64 samples

**Testing**

- no interaction: 36 samples

- meeting: 10 samples

- moving: 16 samples

The split between the unique persons is visualized below:

**ALL UNIQUE PAIRS**

| Pair1 | Pair2 | Pair3 | Pair4 | Pair5 | Pair6 | Pair7 | Pair8 | Pair9 | Pair10 |

**FOLD1**

No interaction:144
Meeting:40
Moving:64

No interaction:36
Meeting:10
Moving:16

**FOLD2**

No interaction:144
Meeting:40
Moving:64

No interaction:36
Meeting:10
Moving:16

**FOLD3**

No interaction:144
Meeting:40
Moving:64

No interaction:36
Meeting:10
Moving:16

Figure 31: 3-fold between distinct pairs of persons

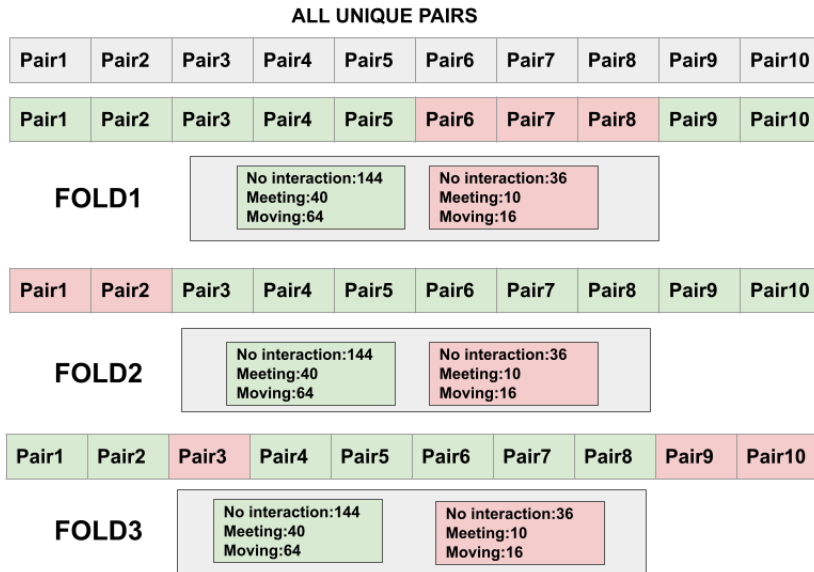As mentioned we have **39 unique pairs** and we have created **3 folds** that each one has a training and testing dataset of the following distribution(with minor differences in the distribution of the classes of the testing dataset):



Figure 32: Complex event distribution in each fold

## 5.4 NeurAsp Learning in Activity Recognition

The deep learning network type chosen for the Complex Event Recognition task is a **Many-to-many Bi-directional LSTM** where a sequence of **T** consecutive frame features are given as input and the output is a matrix of size **1xT** which contains classification results **for each time step of the sequence**.The classification on all time steps is achieved by applying a **softmax activation function** to the last **hidden layer**.

### 5.4.1 Window Size Selection

One very **important parameter** of our task is the **window size** which defines the size of the sequence that needs to be classified. Larger windows require slower training and inference times whereas smaller windows allow faster activity inference and also have reduced needs in terms of resources[Banos, 2014]. There are no best practices for the selection of an optimal window size, and it really depends on the specific model being used and the ideal size is determined through experimentation. We observed that lager window sizes performed better regarding the **discrimination between the classes** of our task. Given that we will try to give a intuitive interpretation why lager windows perform better by visualizing the sequence for three different window sizes. In the images bellow we have drawn the **bounding boxes** and the **center** of a person on a black image so we can observe how the coordinates of the bounding box move through the **consecutive frames**:

- **6 frames**:



active          walking          inactive

Figure 33: Visualized simple events in 6-framed time windows

- **12 frames**:

Figure 34: Visualized simple events in 12-framed time windows

- **24 frames**:


Figure 35: Visualized simple events in 24-framed time windows

As we can observe the more we **increase** the window size the more the classes are **distinctable** especially for the **active** and **walking** classes since the **inactive** class is the most trivial class to identify because the person's bounding box and center **does not change through time**. Given that we have concluded that the window size for the presented experiments of this thesis will be 24.

### 5.4.2 NeurAsp program of the Task

The NeurAsp program for the task is:

```
window(P1). window(P2).
nn(simpleEvent(T,S), [active, inactive, running, walking]) :- window(S).
```

The NeurAsp program expects two input windows of size T (one for each person P1, P2) and will output classifications for each step of the window using the 4 provided labels for the simple events. NeurAsp will internally translate the NeurAsp program to actual ASP code (choice rules) just as we explained in the digit addition example and will generate atoms with all the possible choices for each person in every frame in the form of (for convienience **instead of simpleEvent we write se**):

```
se(0,P1,active), se(0,P1,inactive),se(0,P1,running), se(0,P1,walking),
se(1,P1,active), se(1,P1,inactive),se(1,P1,running), se(1,P1,walking),
```

se(2,P1,active), se(2,P1,inactive),se(2,P1,running), se(2,P1,walking),
se(0,P2,active), se(0,P2,inactive),se(0,P2,running), se(0,P2,walking),
se(1,P2,active), se(1,P2,inactive),se(1,P2,running), se(1,P2,walking),
se(2,P2,active), se(2,P2,inactive),se(2,P2,running), se(2,P2,walking)

For our example of 3 timepoints we have 24 generated atoms.

### 5.4.3 Learning Procedure

The training of our bi-directional LSTM through the NeurAsp network requires the 3rd set of data structures and tensors as described in section **4.3.5** where the input tensors come as pair within a tuple accompanied with their pre-computed stable models.
The training procedure is similar with the digit addition but more complex since we have classifications in each timestep (frame). Let's assume that we have the following tensor pair **of window size T=3** for the **two involved persons p1 and p2**:

| P1 3x5 Tensor | | | | | |
|---|---|---|---|---|---|
| timestep | orientation | xc | yc | height | width |
| 0 | 131 | 184 | 47 | 32 | 34 |
| 1 | 131 | 184 | 47 | 32 | 34 |
| 2 | 131 | 184 | 47 | 32 | 34 |

| P2 3x5 Tensor | | | | | |
|---|---|---|---|---|---|
| timestep | orientation | xc | yc | height | width |
| 0 | 135 | 182 | 45 | 32 | 34 |
| 1 | 135 | 182 | 45 | 32 | 34 |
| 2 | 135 | 182 | 46 | 33 | 34 |

Figure 36: Example of input tensors

The output of the forward pass for the two input tensors is a 3x4 tensor which contains a probability distribution over simple events **(active, inactive, running, walking)** as such:

| P1 3x4 Output tensor | | | | |
|---|---|---|---|---|
| timestep | active | inactive | running | walking |
| 0 | 0.2 | 0.4 | 0.1 | 0.3 |
| 1 | 0.15 | 0.4 | 0.15 | 0.3 |
| 2 | 0.2 | 0.4 | 0.1 | 0.35 |

| P2 3x4 Output tensor | | | | |
|---|---|---|---|---|
| timestep | active | inactive | running | walking |
| 0 | 0.4 | 0.1 | 0.2 | 0.3 |
| 1 | 0.4 | 0.15 | 0.15 | 0.3 |
| 2 | 0.45 | 0.1 | 0.2 | 0.25 |

Figure 36: Example of output tensors

**Note** that the timestep column and the column names in all tensors are used only for illustration purposes. The indexing between the atoms and the cells of a tensor is achieved using the row and the column numbering. Specifically the column indices are the encoded numbers that correspond to the labels 0 ,1, 2, 3 respectively. For instance , row=0 and column=3 refers to the first timestep of the walking simple event probability.

The pre-computed stable models (we assume 3 unique stable models) in the given time range and their probabilities:

**model1**:[ se(0,P1,active), se(0,P2,inactive),se(1,P1,active), se(1,P2,inactive) se(2,P1,active), se(2,P2,inactive) ]

**P(model1)** = P(se(0,P1,active))+P(se(0,P2,inactive))
+ P(se(1,P1,active))+P(se(1,P2,inactive))
+ P(se(2,P1,active))+P(se(2,P2,inactive))
= 0.2×0.4×0.15×0.15×0.2×0.1 = 0.000036

**model2**:[ se(0,P1,inactive), se(0,P2,inactive),se(1,P1,inactive), se(1,P2,inactive) se(2,P1,active), se(2,P2,active) ]

**P(model2)** = P(se(0,P1,inactive))+P(se(0,P2,inactive))
+ P(se(1,P1,inactive))+P(se(1,P2,inactive))
+ P(se(2,P1,active))+P(se(2,P2,active))
= 0.4×0.1×0.4×0.15×0.2×0.45 = 0,000216

**model3**:[ se(0,P1,inactive), se(0,P2,active),se(1,P1,inactive), se(1,P2,active)

se(2,P1,active), se(2,P2,active) ]

**P(model3)** = P(se(0,P1,inactive))+P(se(0,P2,active))
　　　　+ P(se(1,P1,inactive))+P(se(1,P2,active))
　　　　+ P(se(2,P1,active))+P(se(2,P2,active))
　　　　= 0.4×0.4×0.15×0.4×0.2×0.45 = 0,000864

**Note** that the probability of each atom **se(frame_num,person_id,event_label)** map directly to the two output tensors of the neural network as such:

- **person_id** denotes which tensor to lookup(P1 or P2 tensor)

- **frame_num** which row

- **event_label** which simple event column

For instance **P(se(2,P2,inactive)** is referring the to **P2 tensor** , **3rd row** and **simple event's inactive column**.

In a similar way as the digit addition example NeurAsp also keeps a tensor for the gradient updates initialized to zeroes. For our example with T=3 and 4 classes its a 3x4 tensor(the timestep column is used only for representation):

| Gradient Update Tensor | | | | |
|---|---|---|---|---|
| timestep | active | inactive | running | walking |
| 0 | 0 | 0 | 0 | 0 |
| 1 | 0 | 0 | 0 | 0 |
| 2 | 0 | 0 | 0 | 0 |

Figure 36: Gradient tensor filled with zeroes

The next step is the gradient update calculation which is the same iterative procedure with the digit addition example. NeurAsp will iterate through all the atoms created by the choice rules and will compare them against all the stable models one by one. Let's remember the formula of the **NeurAsp loss function:**

57

Figure 37: Formula break down

As we did with digit addition example we will go through a step by step execution for some of the atoms generated by the choice rules.

The Denominator is the sum of probabilities of each stable model that we have already calculated.

$$\sum_{I:|=O} P_{\Pi(\theta)}(I)) = \text{P(model1) + P(model2) +P(model3)}$$
$$= 0,000036 + 0,000216 + 0,000864$$
$$= 0,001116$$

The Numerator and its parts is calculated through the mentioned iterative procedure:

- se(0,P1,active) **is in model1** (part1 of the numerator):

[ se(0,P1,active), se(0,P2,inactive),se(1,P1,active), se(1,P2,inactive), se(2,P1,active), se(2,P2,inactive) ]

There will be a positive update of $\frac{P(model1)}{P(se(0,P1,active))} = \frac{0.000036}{0.2} = 0,0018$

- se(0,P1,active) **is not in model2** (part2 of the numerator):

[ se(0,P1,inactive), se(0,P2,inactive),se(1,P1,inactive), se(1,P2,inactive), se(2,P1,active), se(2,P2,active) ]

There will be a negative update of $\frac{P(model2)}{P(se(0,P1,inactive))} = \frac{0.000216}{0.4} = 0,00054$

- se(0,P1,active) **is not in model3** (part2 of the numerator):

[ se(0,P1,inactive), se(0,P2,active),se(1,P1,inactive), se(1,P2,active),

58

se(2,P1,active), se(2,P2,active) ]

There will be a negative update of $\frac{P(model3)}{P(se(0,P1,inactive))} = \frac{0,000864}{0.4} = 0,00216$

For this execution so far the gradient update from **person P1** for the **weight** of **active simple event** in the **first frame** can be computed:

$$\text{P1\_update} = \frac{0,000036 - 0,00054 - 0,00216}{0,001116} = -2,387$$

The procedure will continue for the rest of the atoms that are referring to P1 and their contribution in terms of gradient updates will be calculated and stored to their respective cells of the gradient tensor. We remind that there is a mapping between the indices of the gradient tensor and the atoms using the **rows to map the timesteps** and **the columns to map the simple event labels**. In the image below we can see the first populated cell while the rest of them are calculated using the procedure we described step by step:

| Gradient Update Tensor | | | | |
|---|---|---|---|---|
| timestep | active | inactive | running | walking |
| 0 | -2.387 | P1_update | P1_update | P1_update |
| 1 | P1_update | P1_update | P1_update | P1_update |
| 2 | P1_update | P1_update | P1_update | P1_update |

Figure 37: Gradient tensor containing the computed gradients of one person

When the iteration reaches the atoms that are referring to P2 the update contributions will be added to the already P1 calculated values. In the image below we can see how the populated gradient matrix looks like after all atoms are checked:

| Gradient Update Tensor | | | | |
|---|---|---|---|---|
| timestep | active | inactive | running | walking |
| 0 | -2.387<br>+<br>P2_update | P1_update<br>+<br>P2_update | P1_update<br>+<br>P2_update | P1_update<br>+<br>P2_update |
| 1 | P1_update<br>+<br>P2_update | P1_update<br>+<br>P2_update | P1_update<br>+<br>P2_update | P1_update<br>+<br>P2_update |
| 2 | P1_update<br>+<br>P2_update | P1_update<br>+<br>P2_update | P1_update<br>+<br>P2_update | P1_update<br>+<br>P2_update |

Figure 37: Gradient tensor containing the computed gradients of both persons

Finally the gradient update tensor is computed, gets multiplied by the value of **learning_rate** and when the number of samples that are contributing to gradient updates reaches the value of **batch_size**, the **backpropagation is triggered**.

## 5.5   NeurAsp Inference in Activity Recognition

In this section we will demonstrate how the outputs of the LSTM are integrated with ASP code. We will go through on how simple event atoms are generated using the LSTM outputs and how those atoms are used to infer on complex events by combining the rules we introduced in section **4.3.4 Ground Truth Generation on Complex Events with ASP**.

### 5.5.1   Simple Events atoms generation from LSTM outputs

Since the output tensors of the neural network is a probability distribution over simple events we apply an **argmax operation** to retrieve the event with the highest probability as shown in the picture below:

| P1 3x4 Output tensor | | | | |
|---|---|---|---|---|
| timestep | active | inactive | running | walking |
| 0 | 0.2 | 0.4 | 0.1 | 0.3 |
| 1 | 0.15 | 0.4 | 0.15 | 0.3 |
| 2 | 0.2 | 0.4 | 0.1 | 0.35 |

MAX →

| P1 simple event result | |
|---|---|
| timestep | result |
| 0 | inactive |
| 1 | inactive |
| 2 | inactive |

| P2 3x4 Output tensor | | | | |
|---|---|---|---|---|
| timestep | active | inactive | running | walking |
| 0 | 0.4 | 0.1 | 0.2 | 0.3 |
| 1 | 0.4 | 0.15 | 0.15 | 0.3 |
| 2 | 0.45 | 0.1 | 0.2 | 0.25 |

MAX →

| P2 simple event result | |
|---|---|
| timestep | result |
| 0 | active |
| 1 | active |
| 2 | active |

Figure 38: Predictions for each person

Using the mapping of each output tensor we can generate asp atoms in the form: **simpleEvent(frame_num,person_id,simple_ev_label).** For instance in our depicted example we will retrieve the following simple event atoms:

> **P1 atoms**: se(0,P1,inactive), se(1,P1,inaswsctive), se(2,P1,inactive)
> **P2 atoms**: se(0,P2,active), se(1,P2,active), se(2,P2,active)

This way we have encoded the LSTM outputs as ASP atoms. We remind that the that each sample consists of two persons involved and the classification outputs of the LSTM comes one at the time so first we obtain the output results of P1 and then the results of P2.

### 5.5.2   Application of logic on LSTM outputs with ASP

We remind that we will infer on three complex event classes: **no_interaction, meeting, moving**. The **no_iteraction** complex event holds when neither of the other two holds. Regarding the **meeting** and **moving** complex events we have defined them as ASP code where the law of inertia was used as a basis. Besides the atoms generated from LSTM outputs the rules need additional information in order to operate such as the **coordinates , visibility and orientation atoms**. This additional information comes alongside with each time window sample as we have described in **4.3.5 Input and Label Tensors** where we keep a **list of atoms** extracted directly from the XML files. To summarize we need the following to perform inference using the NeurAsp framework:

- **atoms from LSTM outputs**

- **information atoms directly from XMLs**

- **complex event rules definitions.**

For the ASP part that we have described so far we can imagine an additional **"logic layer"** at the end of our LSTM where the complex event definitions reside as **background knowledge (BK)** , takes as inputs the **P1 and P2 LSTM output atoms** and the **additional XML information atoms**. By using **clingo** the logic layer outputs results on complex events. An illustration is given below:



Figure 38: Diagram on how the logic layer fits

# 6 Experimental Methods

In this section we will describe the methods used in order to investigate how the **application of logic** could benefit deep learning methods in terms of **inference** and **training** in the Complex Event Recognition Task.

We remind that the **complex events** that we want to classify are 3: **no interaction, meeting, moving** whereas the simple events are 4: **active, inactive, walking, running**.

In the data pre-processing procedure we have created **3 sets of of data structures and tensors** which correspond to **3 different approaches**:

- **Only deep learning methods/method1**: In this setup we will train a neural network to **classify the complex events** by giving as input a single sequence tensor that contains features for two persons.

- **Traditional neural network training and inference with logic/method2**: In this setup we train a neural network to **classify simple events** by giving as input individual sequence tensors that contain features for one person at the time. The **ground truth for the simple events is provided**

**directly from the CAVIAR dataset.**The **classification on complex events** is achieved by applying ASP rules at the output layer.

- **NeurAsp Integration/method3**: Finally we will use the **NeurAsp framework** for training and inference. In this setup **no ground truth will be provided for the simple events but rather stable ASP models** that contain simple events. Again the **classification on complex events** is achieved by applying ASP rules at the output layer

## 6.1 Bi-directional LSTM Hyperparameters

The bi-directional LSTM is the **core component** of all 3 approaches and in order to make a concrete **comparison** we have chosen the same **main hyperparameters** for all methods. There are differences between the some hyperparameters:

- **input layer**: in method1 we provided a tensor that contains the concatenated features of the two persons.

- **output layer**: in method1 the output of LSTM contains predictions on complex events whereas in the two remaining methods the output contains predictions on simple events.

- **loss function**: in method3 we do not use the cross entropy loss since we use the **semantic loss** described in **Learning with NeurAsp section**.

Finally we provide the complete list of the overall hyperparameters:

- **Loss function**: crossEntropy (method1,method2) / semantic (method3)

- **Learning rate**: 0.0000015

- **Number of hidden layers**: 2

- **Number of neurons in each hidden layer**: 256

- **Dropout**: 0.3

- **Epochs**: 1000

- **Batch size**: 8

- **Input Layer size**: 11 (method1) / 5 (method2 , method3)

- **Softmax Output Layer size**: 3 (method1) / 4 (method2 , method3)

- **window size T**: 24

The stratified 3-fold split has yielded 3 folds where each fold is split to training and testing/unseen data. For every method we let the LSTM train for 1000 epochs for each fold. During the training procedure we keep the best macro f1 score on the training data. After the training procedure is completed for each fold the results are compared to unseen data as selected from the stratified 3-fold validation.

For every method we provide the results of each fold individually by providing the confusion matrices and the tables that contain more detailed evaluation metrics(precision , recall , f1-score).The macro f1-score of each fold is computed by averaging the individual f1-scores for each complex event class.

Finally we aggregate the **macro average f1-scores** of each fold to form the final evaluation metric of the method. The macro average f1-score as we have mentioned is the most suitable evaluation metric for imbalanced datasets such as ours.

In the first method where no logic is applied the LSTM outputs results **directly on complex events**. On the other hand the two remaining methods use the **"logic layer"** we described in the previous section.

Finally it is important to emphasize that **the essential difference** between the three methods besides the **inference method** is the **computation of the loss and the gradient updates** so we will provide an illustration of the training pipeline for each method.

## 6.2   Method1: Only Deep Learning methods

This method is the most straightforward since it does not involve any logic integration.The input tensor is of size 24x11 which **contains the concatenated features of both persons plus the distance between their coordinates**.

### 6.2.1   Method1 Pipeline Details

This method uses a traditional deep learning pipeline. In the picture below we can see two distinct phases of the pipeline, the **training phase** and the **inference phase**:

Figure 39: Method 1 diagram

- **training phase**: The tensor with the concatenated features is given as input to the LSTM and an output tensor is produced which contains a probability distribution over complex events for each time step. This output tensor is propagated along with the ground truth tensor to the Cross Entropy Loss. Then using the **gradient descent algorithm** the gradient updates that minimize the loss are computed and **backpropagated** to the network when the number of processed samples reaches the **batch size**.

- **inference phase**: The inference phase is achieved by applying an **argmax** operation to the output tensor and obtain the classes that have the highest probability. The network is evaluated by comparing the predictions with the ground truth.

### 6.2.2 Method1 Results

## Fold1 Results

| Confusion Matrix | Evaluation Metrics Tables |
|---|---|

| | no_int | meeting | moving |
|---|---|---|---|
| no_int | 1712 | 96 | 232 |
| meeting | 86 | 91 | 241 |
| moving | 809 | 0 | 429 |

| Complex Event | Precision | Recall | f1-score |
|---|---|---|---|
| no_int | 0.66 | 0.84 | 0.74 |
| meeting | 0.49 | 0.22 | 0.30 |
| moving | 0.48 | 0.35 | 0.40 |
| **Overall Metrics** | | | |
| Accuracy | 0.60 | | |
| Macro f1-score | 0.48 | | |

## Fold2 Results

| Confusion Matrix | Evaluation Metrics Tables |
|---|---|

| | no_int | meeting | moving |
|---|---|---|---|
| no_int | 1220 | 414 | 406 |
| meeting | 36 | 295 | 82 |
| moving | 685 | 16 | 542 |

| Complex Event | Precision | Recall | f1-score |
|---|---|---|---|
| no_int | 0.63 | 0.60 | 0.61 |
| meeting | 0.41 | 0.71 | 0.52 |
| moving | 0.53 | 0.44 | 0.48 |
| **Overall Metrics** | | | |
| Accuracy | 0.56 | | |
| Macro f1-score | 0.54 | | |

## Fold3 Results

| Confusion Matrix | Evaluation Metrics Tables |
|---|---|

| | no_int | meeting | moving |
|---|---|---|---|
| no_int | 1718 | 123 | 199 |
| meeting | 29 | 311 | 76 |
| moving | 730 | 0 | 508 |

| Complex Event | Precision | Recall | f1-score |
|---|---|---|---|
| no_int | 0.69 | 0.84 | 0.76 |
| meeting | 0.72 | 0.75 | 0.73 |
| moving | 0.65 | 0.41 | 0.50 |
| **Overall Metrics** | | | |
| Accuracy | 0.69 | | |
| Macro f1-score | 0.67 | | |

**Average macro f1 score = 0.56**

## 6.3 Method2: Training on simple events and Infer with NeurAsp

In this method we use the **ground truth of simple events as provided directly from the annotators in the XML files**. The input is a pair of tensors of size 24x5 that contain the features of each person individually.

### 6.3.1 Method2 Pipeline Details

This method for the training phase uses the traditional deep learning approach and uses the "logic layer" to produce predictions:



Figure 40: Method 2 diagram

- **training phase**: Here we have two individual tensors that are given one at the time. When the LSTM produces their respective output tensors they are sent alongside with their ground truth on **simple events** to the cross entropy loss function. Similarly with the first method the updates are computed using the **gradient descent algorithm** and then back-propagation is initiated when the number of samples reaches the **batch size**. (the batch size for this method is at minimum 2 since the tensors come as a pair so in order to have the same batch size as the previous method we declare a batch size of 16)

- **inference phase**: After the LSTM has produced an output tensor for each person the **argmax** operation is applied individually and we get the predictions on simple events. Those predictions are translated to **simple events atoms** (se atoms) and alongside with the needed information XML

67

atoms (coordinates, visibility ,etc) are given to the **"logic layer"**. The "logic layer" outputs predictions on complex events using **clingo** and the evaluation of the network is done by comparing the predictions with the ground truth on complex events.

### 6.3.2 Method2 Results

#### Fold1 Results

| Confusion Matrix | Evaluation Metrics Tables | | |
|---|---|---|---|
| | **Complex Event** | **Precision** | **Recall** | **f1-score** |

| Complex Event | Precision | Recall | f1-score |
|---|---|---|---|
| no_int | 1 | 1 | 1 |
| meeting | 1 | 0.13 | 0.23 |
| moving | 0.92 | 90 | 0.87 |

| Overall Metrics | |
|---|---|
| Accuracy | 0.92 |
| Macro f1-score | 0.79 |

Confusion Matrix (rows no_int, meeting, moving; columns no_int, meeting, moving):

|  | no_int | meeting | moving |
|---|---|---|---|
| no_int | 2040 | 0 | 0 |
| meeting | 0 | 55 | 363 |
| moving | 0 | 0 | 1238 |

#### Fold2 Results

| Complex Event | Precision | Recall | f1-score |
|---|---|---|---|
| no_int | 0.95 | 0.96 | 0.96 |
| meeting | 1 | 0.55 | 0.71 |
| moving | 0.85 | 0.97 | 0.91 |

| Overall Metrics | |
|---|---|
| Accuracy | 0.92 |
| Macro f1-score | 0.86 |

Confusion Matrix (rows no_int, meeting, moving; columns no_int, meeting, moving):

|  | no_int | meeting | moving |
|---|---|---|---|
| no_int | 1951 | 0 | 89 |
| meeting | 58 | 227 | 128 |
| moving | 35 | 0 | 1208 |

**Fold3 Results**

| Confusion Matrix | Evaluation Metrics Tables | | | |
|---|---|---|---|---|

Confusion Matrix:

|  | no_int | meeting | moving |
|---|---|---|---|
| **no_int** | 1954 | 0 | 86 |
| **meeting** | 18 | 269 | 131 |
| **moving** | 0 | 0 | 1238 |

| Complex Event | Precision | Recall | f1-score |
|---|---|---|---|
| no_int | 0.99 | 0.96 | 0.97 |
| meeting | 1 | 0.64 | 0.78 |
| moving | 0.85 | 1 | 0.92 |

| Overall Metrics | |
|---|---|
| Accuracy | 0.94 |
| Macro f1-score | 0.89 |

**Average macro f1 score = 0.82**

## 6.4  Method3: NeurAsp Training and Inference

This final method uses the NeurAsp training for both **training** and **inference**. The LSTM is trained to classify simple events just like method2 but **instead of providing ground truths** we provide **stable models** that we generated using the procedure we described in subsection **Stable Model generation for time windows**.

### 6.4.1  Method3 Pipeline Details

The method3 pipeline seems quite similar with the pipeline of the previous method since the inference phase is the same but there is a big difference in the training procedure:

Figure 41: Method 3 diagram

- **training phase**: Again we have two individual tensors that are given one at the time. As we have described in section **NeurAsp Learning on Complex Event Recognition** the NeurAsp program has produced asp atoms using choice rules for all possible simple event outcome for every timestep. Those atoms are mapped with the probabilities that the LSTM has produced and the iterative procedure that checks each stable model begins. The backpropagation is initiated when the number of processed samples is equal to the batch size.

- **inference phase**: The inference method is exactly the same as method2.

### 6.4.2  Method3 Results

## Fold1 Results

| Confusion Matrix | Evaluation Metrics Tables | | |
|---|---|---|---|
| | **Complex Event** | **Precision** | **Recall** | **f1-score** |

**Confusion Matrix (Fold1):**

|  | no_int | meeting | moving |
|---|---|---|---|
| **no_int** | 2040 | 0 | 0 |
| **meeting** | 113 | 240 | 65 |
| **moving** | 74 | 22 | 1142 |

| Complex Event | Precision | Recall | f1-score |
|---|---|---|---|
| no_int | 0.95 | 1 | 0.97 |
| meeting | 0.94 | 0.41 | 0.57 |
| moving | 0.87 | 0.96 | 0.92 |

| Overall Metrics | |
|---|---|
| Accuracy | 0.92 |
| Macro f1-score | 0.82 |

## Fold2 Results

**Confusion Matrix (Fold2):**

|  | no_int | meeting | moving |
|---|---|---|---|
| **no_int** | 1957 | 0 | 83 |
| **meeting** | 48 | 216 | 149 |
| **moving** | 0 | 0 | 1243 |

| Complex Event | Precision | Recall | f1-score |
|---|---|---|---|
| no_int | 0.98 | 0.96 | 0.97 |
| meeting | 1 | 0.52 | 0.69 |
| moving | 0.84 | 1 | 0.91 |

| Overall Metrics | |
|---|---|
| Accuracy | 0.92 |
| Macro f1-score | 0.86 |

**Fold3 Results**

| Confusion Matrix | | | Evaluation Metrics Tables | | | |
|---|---|---|---|---|---|---|

| | Complex Event | Precision | Recall | f1-score |
|---|---|---|---|---|
| no_int | 2040 | 0 | 0 |
| meeting | 24 | 336 | 58 |
| moving | 59 | 4 | 1175 |
| | no_int | meeting | moving |

| Complex Event | Precision | Recall | f1-score |
|---|---|---|---|
| no_int | 0.96 | 1 | 0.98 |
| meeting | 0.99 | 0.80 | 0.89 |
| moving | 0.95 | 0.95 | 0.95 |
| **Overall Metrics** | | | |
| Accuracy | 0.96 | | |
| Macro f1-score | 0.94 | | |

**Average macro f1 score = 0.87**

# 7  Experimental Evaluation

In this section we will summarize and compare the results of the presented methods. The stratified 3-fold has provided three sets of training and testing datasets where each fold consists of 3344 training sample frames (or 556 time windows of size 24) and 3696 testing sample frames (or 154 time windows of size 24). All methods are trained for 1000 epochs for each fold. The table below compares the **average macro f1 score** calculated from all folds and the training time needed from each method for one fold.

| Method | Avg Macro f1 score | Training Time |
|---|---|---|
| Only Neural | 0.56 | 37 mins |
| Neural Train and NeurAsp Inf | 0.82 | 47 mins |
| NeurAsp Train and Inf | 0.87 | 62 mins |

Table 5: Overall results

The **first method** (**Only Neural**) has by far the worst performance where the input was the concatenated features of the two involved individuals and the classification was performed directly on complex events. In fact the LSTM could not make a clear distinction between the no interaction class and the two remaining classes.

Following the **second method** (**Neural Train and NeurAsp Inf**) there is a significant improvement regarding the performance. The addition of the logic layer at the end of the LSTM proves that the integration of neural and logic methods could yield better results in this specific task.

Finally the **third method** (**NeurAsp Train and Inf**) has the best results

of all the presented methods. Using the NeurAsp framework for training has added an overhead in training time due to the internal iterative procedure but the results show a clear improvement regarding the overall performance.

# 8    Results Summary  Discussion

The aim of this thesis was to investigate and present the capabilities of neuro-symbolic AI through the usage of NeurAsp framework. We have presented how it integrates deep-learning and logic-based methods, how it works internally and finally we applied the challenging task of human activity recognition in videos. With the presentation of three different experimental setups which compared a pure neural approach, a method that apply logic mechanisms in inference and finally a method that uses logic both in inference and training, we have concluded that neuro-symbolic methods could yield promising results. Given that CAVIAR is an old dataset with low resolution images we used the provided features that we mentioned such as coordinates and bounding box sizes. In case we had a dataset with good resolution images we could extract more informative features that involve human body, hand, facial, and foot keypoints. There are several frameworks and systems with impressive results in real time video such as openpose[4] and and media pipe pose[5]. Many approaches in HAR use those features [Sawant, 2020] [Noori, 2019] so a future work could be focused on creating logic rules that apply on body keypoints. Finally since the position of the camera or actors is one of the biggest challenges in HAR, future works could focus on creating logic rules dedicated to the camera angle and the overall homography of the area of interest.

# References

[1] Zachary Susskind, Bryce Arden, Lizy K. John, Patrick Stockton, Eugene B. John Neuro-Symbolic AI: An Emerging Class of AI Workloads and their Characterization

[2] Rocktäschel et al, End-to-End Differentiable Proving. NeurIPS, 2017

[3] Sourek et al, Lifted relational neural networks: Efficient Learning of Latent Relational Structures, JAIR 2018

[4] Evans et al, Learning Explanatory Rules from Noisy Data, JAIR, 2018

[5] Diligenti et al, Semantic-Based Regularization for Learning and Inference, AIJ, 2017

[6] Xu et al, A Semantic Loss Function for Deep Learning with Symbolic Knowledge, ICML, 2018

---

[4]https://github.com/CMU-Perceptual-Computing-Lab/openpose
[5]https://google.github.io/mediapipe/solutions/pose.html

[7] Marra et al, Neural Markov Logic Networks, UAI, 2021

[8] Badreddine et al, Logic Tensor Networks, AIJ, 2022

[9] Robin Manhaeve, Sebastijan Du-mancic, Angelika Kimmig, Thomas De-meester, and LucDe Raedt. Deepproblog: Neural probabilistic logic pro-gramming. In Proceedings of Advances in Neural Infor-mation Processing Systems, pages 3749–3759, 2018

[10] Yang et al, NeurAsp: Embracing Neural Networks Into Answer Set Pro-gramming, IJCAI 2020.

[11] Winters et al, DeepStochLog: Neural stochastic logic programming, AAAI, 2022

[12] Garcez A. et al, Neural-symbolic computing: An Effective Methodology for Principled Integration of Machine Learning and Reasoning. FLAP, 2019.

[13] Lifschitz, V. Answer set programming. Springer, 2019.

[14] Neha Gupta, Suneet K Gupta, Rajesh K Pathak, Vanita Jain, Parisa Rashidi, Jasjit S Suri Human activity recognition in artificial intelligence framework: a narrative review

[15] Gulustan Dogan, Sinem Sena Ertas, İremnaz Cay Human Activity Recog-nition Using Convolutional Neural Networks

[16] Deepika Singh, Erinc Merdivan, Ismini Psychoula, Johannes Kropf, Sten Hanke, Matthieu Geist, Andreas HolzingerHuman Activity Recognition us-ing Recurrent Neural Networks

[17] Keze Wang, Xiaolong Wang, Liang Lin, Meng Wang, Wangmeng Zuo 3D Human Activity Recognition with Reconfigurable Convolutional Neural Net-works

[18] Michalis Vrigkas, Christophoros Nikou, Ioannis A. Kakadiaris A Review of Human Activity Recognition Methods

[19] Nikos Katzouris Scalable Relational Learning for Event Recognition

[20] Cugola, G. and Margara, A. (2012). Processing flows of information: From data stream

[21] Nikos Giatrakos, Elias Alevizos, Alexander Artikis,Antonios Deligiannakis, Minos Garofalakis Complex event recognition in the Big Data era: a survey

[22] Alexander Artikis, Marek Sergot, Georgios Paliouras An Event Calculus for Event Recognition

[23] Elias Alevizos Complex Event Forecasting A Formal Framework , Section 2.3

[24] Katzouris,Paliouras, Artikis Online Learning Probabilistic Event Calculus Theories in Answer Set Programming

[25] Javier Ortiz Laguna, Angel García Olaya, Daniel Borrajo A Dynamic Sliding Window Approach for Activity Recognition

[26] Nikos Katzouris, Alexander Artikis WOLED: A Tool for Online Learning Weighted Answer Set Rules for Temporal Reasoning Under Uncertainty https://cer.iit.demokritos.gr/publications/papers/2020/kr-2020.pdf

[27] Oresti Banos,Juan-Manuel Galvez, Miguel Damas, Hector Pomares, Ignacio Rojas Window Size Impact in Human Activity Recognition

[28] Alevizos, E., Skarlatidis, A., Artikis, A. and Paliouras, G. 2017. Probabilistic complex event recognition: A survey. ACM Computing Surveys 50, 5, 71:1–71:31. https://dl.acm.org/doi/10.1145/3117809

[29] Chinmay Sawant Human activity recognition with openpose and Long Short-Term Memory on real time images

[30] Farzan Majeed Noori, Benedikte Wallace, Md. Zia Uddin, Jim Torresen A Robust Human Activity Recognition Approach Using OpenPose, Motion Features, and Deep Recurrent Neural Network