
OmniNet: A unified architecture for multi-modal multi-task learning

Subhojeet Pramanik *
IBM Watson
email@subho.in

Priyanka Agrawal
IBM Research
pagrawal.ml@gmail.com

Aman Hussain
University of Amsterdam
email@amanhussain.com

Abstract

Transformer is a popularly used neural network architecture, especially for language understanding. We introduce an extended and unified architecture which can be used for tasks involving a variety of modalities like image, text, videos, etc. We propose a spatio-temporal cache mechanism that enables learning spatial dimension of the input in addition to the hidden states corresponding to the temporal input sequence. The proposed architecture further enables a single model to support tasks with multiple input modalities as well as asynchronous multi-task learning, thus we refer to it as *OmniNet*. For example, a single instance of *OmniNet* can concurrently learn to perform the tasks of part-of-speech tagging, image captioning, visual question answering and video activity recognition. We demonstrate that training these four tasks together results in about three times compressed model while retaining the performance in comparison to training them individually. We also show that using this neural network pre-trained on some modalities assists in learning an unseen task. This illustrates the generalization capacity of the self-attention mechanism on the spatio-temporal cache present in *OmniNet*.²

1 Introduction

Transformer [1] is currently one of the best performing models for any sequence transduction tasks, especially those involving natural language. It is originally designed for a single task at a time. In fact, most of the generic deep learning architectures [2, 3, 4] that have been designed and developed are able to learn, albeit very well, a single task and handle one task specific input domain like image, text or audio. Furthermore with these models, we often rely on the generalization capability of the trained network to guarantee performance on unseen examples. Transfer learning [5, 6] is another popular paradigm used to adapt the model to learn a related task with similar input domain. The success of neural networks across these challenges is known to be due to their ability in learning effective representations of the data. For example, the self-attention mechanism in Transformers can capture the global temporal dependence in sequential data very well. Naturally, the question arises whether we can extend these architectures, like the Transformer, to be able to learn shared representations from multiple input domains and to be able attend on them representations to perform a multitude of tasks concurrently.

The research into multi-task models that learn to solve varied tasks across a multitude of input domains is not new. Work done in [7] demonstrates an architecture capable of learning a shared representation across audio and video modalities. Similarly in [8] a convolutional architecture has been designed to support a variety of NLP tasks. However, most of these architectures are designed to learn specific set of tasks with known input domains. To the best of our knowledge, there does not exist a single unified architecture that works out of the box for any combination of multi-modal tasks.

*Corresponding author

²Source code available at: <https://github.com/subho406/OmniNet>

To address this gap, we introduce a unified architecture, namely *OmniNet*, which enables a single model to support tasks with multiple input modalities and asynchronous multi-task learning. *OmniNet* mimics a computer in terms of input peripherals to support any input modality and a Central Neural Processor (CNP) that enables multi-task learning. We consider that most real-life data like image, text, speech, video, etc. is a direct conjunction of spatial and temporal components. Therefore, we employ a spatio-temporal cache mechanism to learn a shared representation of the input data across the spatial (space) and temporal (time) dimension. Using a generalized *encode()* function, *OmniNet* can process and store spatio-temporal representation for each of the input domains and then *decode()* predictions across a multitude of tasks. In our experiments, we train a single compressed instance of the *OmniNet* to solve several multi-domain tasks such as part-of-speech tagging, image captioning, visual question answering and video activity recognition. To make our work reproducible, open to scrutiny and further development, we have open sourced a demonstration of our system implemented using Pytorch [9] at <https://github.com/subho406/OmniNet>.

2 Related Work

Multi-task learning has been extensively studied in the literature, with applications to a wide set of problems ranging from natural language processing (NLP) [8, 10, 11] to speech recognition [12, 13] to vision [14, 15, 16]. It has also found its use in a combination of diverse tasks like image captioning and text translation and parsing [17, 18]. However, most of these architectures assume the set of tasks to be known in advance. Similarly, multi-modal learning has been essential for solving a broad range of interesting problems such as Visual Question Answering [19, 20] and Video Question Answering [21]. Again, the state-of-the-art models are highly specific to the objective in hand and not easily adaptable to different tasks or domains.

The closest work towards a universal multi-modal multi-task model is the MultiModel [22] architecture which demonstrated for the first time, the ability to learn a multitude of tasks of varying input domains. However, the MultiModel architecture lacks support for multi-modal tasks with more than one input domains such as visual question answering. Further, it only attends over the temporal states of the input data and includes no mechanism to attend over the spatial aspects of the data. To the best of our knowledge, our architecture is the first to unify the encoding process on any number of input domains and includes mechanisms to jointly attend on spatial and temporal representations of multiple input domains at the same time.

3 Proposed Model

We propose a unified architecture, namely *OmniNet*, to enable learning multi-modal tasks with multiple input domains and support generic multi-tasking for any set of tasks. The *OmniNet* architecture consists of multiple sub-networks, called peripheral networks, connected to a common central neural network called the Central Neural Processor (CNP) (Figure 1). Each peripheral network is used to encode the domain specific input into feature representations. In this work, we describe image, text and video peripherals (Section 3.1). One can add more, say speech peripheral, depending on the task. The output representation of a peripheral network is always a spatio-temporal tensor $x \in \mathbb{R}^{t \times s \times d_{model}}$, where t & s are the temporal and spatial dimensions of the input respectively, and d_{model} is the model dimension input to the Central Neural Processor.

The spatio-temporal representations generated by the peripheral networks corresponding to each input domain are then processed by the Central Neural Processor (CNP). The CNP uses fully attention based encoder-decoder [23, 24, 25] model for sequence transduction similar to the Transformer architecture [1], which is the state-of-the-art for multiple language modeling tasks (Section 3.2). During the encoding stage, the CNP implements a generic *encode(x, D)* function to first process and store the spatio-temporal representations of the input, where $x \in \mathbb{R}^{t \times s \times d_{model}}$ is the spatio-temporal tensor produced by the peripheral networks and $D \in \mathbb{Z} : 0 \leq D < D_{len}$ is the domain id and D_{len} is the max number of domains supported by the CNP. The *encode()* function is called multiple times, once for each multi-modal input from respective peripheral. During the decoding stage, a *decode($y_{shifted}, \tau$)* function is used to decode predictions as softmax probabilities, where $y_{shifted} \in \mathbb{Z}^{N-1}$ are the target outputs shifted one time-step to the right, N is the length of the output sequence; $\tau \in \mathbb{Z} : 0 \leq \tau < \tau_{len}$ is task id and τ_{len} is the total number of supported tasks. The

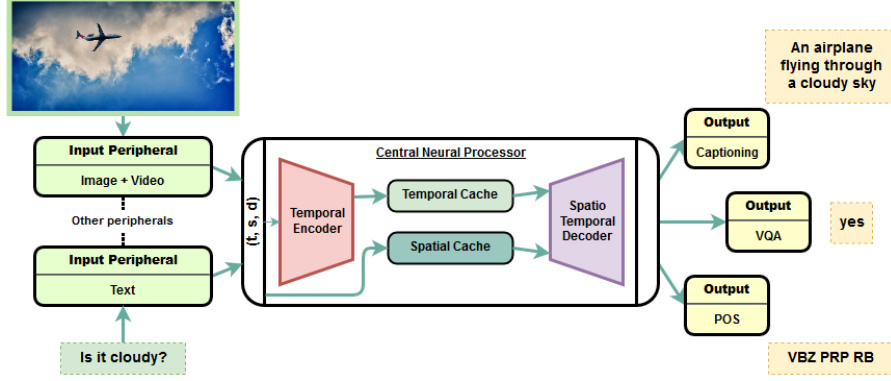


Figure 1: *OmniNet* performing image captioning, visual question answering and POS tagging at once

decoding step is similar to [1], modified to incorporate a two-step attention mechanism over spatial and temporal cache.

3.1 Peripheral networks

First, we elaborate on how we support multiple input domains using peripheral networks. A peripheral network can use a pre-trained model from existing literature to ultimately encode a given domain input to a standardized feature representation $x \in \mathbb{R}^{t \times s \times d_{model}}$, where t & s are the temporal and spatial dimensions of the input respectively, and d_{model} is the model dimension input to the Central Neural Processor. Here we detail text and vision peripherals and one can add more peripherals or alter the peripheral design depending on the task.

Vision peripheral: This peripheral uses a convolutional neural network to encode image and video inputs in the tasks. In our experiments, we use the pre-trained *ResNet-152* model, a variant of ResNet [26] consisting of 152 convolutional layers. We removed the final fully connected and avg-pooling layers to generate spatial feature representations for a given image/video. For an image of dimension $h \times w \times n_c$, *ResNet* down-samples the image into $h' \times w' \times n'_c$, where h, w, n_c are the height, width and number of input channels. For a video, each frame is input to the peripheral to produce $F \times h' \times w' \times n'_c$, where F is the total number of frames in the video. The encoding vectors are then projected to dimension d_{model} using a fully connected layer. The output is then reshaped into a spatio-temporal tensor of $x \in \mathbb{R}^{t \times h' \times w' \times d_{model}}$, where t equal to 1 for an image and F for a video.

Language peripheral: The Language peripheral uses byte-pair encoding [27] to generate subwords for a given input sentence. The subwords are passed to an embedding layer to generate subword embeddings of dimension d_{emb} and projected to dimension d_{model} using a fully connected layer. We used pre-trained subword embeddings with $d_{emb} = 300$ and $vocab_size = 25000$ from [28], which includes pre-trained subword embeddings of over 275 languages, to initialize the weights of the embedding matrix. The output is then reshaped into a spatio-temporal tensor $x \in \mathbb{R}^{t \times 1 \times d_{model}}$, where t equal to number of subwords in the input sentence. As we do not have any spatial dimension in textual data, the spatial dimension of x from a Language peripheral is always 1.

3.2 Central Neural Processor (CNP)

To process the spatio-temporal information in the input data, the CNP implements a spatial cache C_s , temporal cache C_t and a link array L . The spatial and temporal cache and the link array are a list of elements, initialized as empty before the encoding process. During the encoding stage, an `encode()` routine takes as input, the tensor x generated from the peripheral and corresponding domain/peripheral id D . This function processes the spatial and temporal information in the input x and stores them into the spatial cache C_s and the temporal cache C_t , respectively and stores their dimensions t & s in the link array. For a given task, this `encode()` routine is called K times, where K is the number of inputs in the task. Note that these inputs can belong to same or different domains. The function is defined as follows:

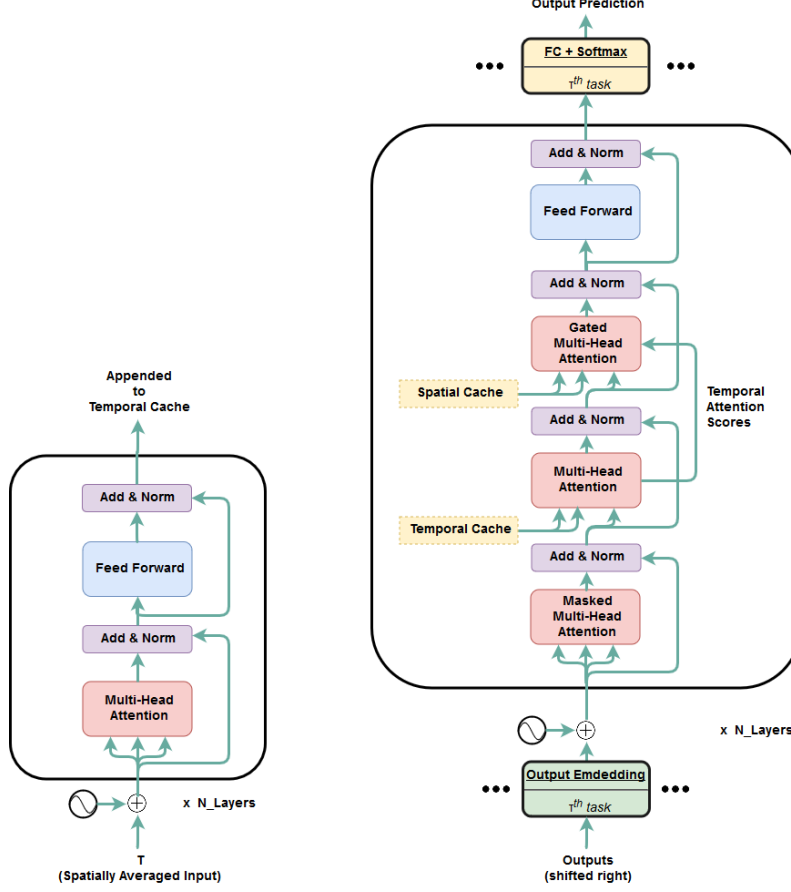


Figure 2: **left:** *TemporalEncoder* architecture; **right:** *OmniNet decode()* architecture.

Encode (x, D): For a given input $x \in \mathbb{R}^{t \times s \times d_{model}}$ and domain identifier D , the *encode()* routine is described in Algorithm 1. Since inputs can come from multiple peripherals, the algorithm first concatenates the input with the domain embedding to ensure a domain-aware encoding of the input (Steps 2 to 3). Steps 4 to 7 process the spatial information in x by unrolling the time dimension and adding these unrolled vectors into the spatial cache. Steps 8 to 10 process the temporal information in x by averaging the spatial dimension of x and then passing the averaged tensor to a self-attention based *TemporalEncoder*. This *TemporalEncoder* is similar to the encoder used in [1] as shown in Figure 2 is used to calculate temporal embeddings of the input sequence. The output from the *TemporalEncoder* is appended to the temporal cache.

The above encoding routine keeps appending spatio-temporal information to C_t & C_s for each input $x_k \in \mathbb{R}^{t_k \times s_k \times d_{model}}$. Note the superscript k to denote correspondence to k -th input of the task, where $k \in 1, \dots, K$. After K calls, we have the temporal cache $C_t = [T_1, \dots, T_R]$, where $R = \sum_{r=1}^K t_r$; the spatial cache $C_s = [S_1, \dots, S_P]$, where $P = \{\sum_p t_p * s_p : p \in 1, \dots, K \wedge s_p > 1\}$ and the link array $L = [(t_1 \rightarrow s_1), \dots, (t_K \rightarrow s_K)]$. Note that C_s can also be empty in case the *encode()* is only called with inputs with $s_k = 1 \forall k$. Next, we use the *decode()* function to generate predictions as softmax probabilities:

Decode ($y_{shifted}, \tau$): The architecture of the *decode()* function is shown in Figure 2. The *decode()* takes as argument the output labels $y_{shifted}$ shifted one time step to the right, a task id τ and generates predictions by attending from the spatial and temporal cache. The *decode()* function is structured similar to the decoder used in the *Transformer* architecture [1] and jointly attends on the vectors stored in the temporal and spatial cache. Similar to [1], the decoding first starts by attending over the output embeddings using masked multi-head scaled dot product attention. The attention layer for the temporal cache uses scaled dot-product attention with multiple heads as specified in [1]. The attention

Algorithm 1 *encode()*: Encodes spatial and temporal representations into spatial and temporal cache

Require: $x \in \mathbb{R}^{t \times s \times d_{model}}$, D , C_s , L , C_t

```

1:  $L \leftarrow L \cup (t \rightarrow s)$ 
2:  $D_{emb} \leftarrow EmbedLayer(D)$ 
3:  $x \leftarrow FC(Concat(x, D_{emb}), d_{model})$ 
4: if  $s > 1$  then
5:    $S \leftarrow Reshape(x, (ts, d_{model}))$  {where, output  $S = [S_1, \dots, S_{ts}]$  s.t.  $S_i \in \mathbb{R}^{d_{model}}$  is a
     spatial feature vector.}
6:    $C_s \leftarrow C_s \cup [S_1, \dots, S_{ts}]$  {Append spatial representations to spatial cache}
7: end if
8:  $T \leftarrow (\sum_{i=1}^s x[:, i, :]) / s$ 
9:  $T \leftarrow TemporalEncoder(T)$  {where, output  $T = [T_1, \dots, T_t]$  s.t.  $T_j \in \mathbb{R}^{d_{model}}$  is the encoding
   of temporal dimension in  $x$ .}
10:  $C_t \leftarrow C_t \cup [T_1, \dots, T_t]$  {Append temporal representations to temporal cache}

```

layer for the spatial cache, uses gated multi-head attention to attend over the elements of the spatial cache. For inputs with both time and space dimension (e.g. video), we want the spatial attention layer to attend more on frames which have relatively high attention scores in the temporal cache attention layer. Therefore, the attention score output from the temporal cache multi-head attention layer $A \in \mathbb{R}^{n_h \times N \times R}$, is used to calculate the tensor $G \in \mathbb{R}^{n_h \times N \times P}$ used for gating the attention score output in the spatial attention layer, where n_h is the number of heads in multi-head attention as described in [1]. The tensor G is calculated using A & L as detailed in Algorithm 2.

Algorithm 2 Calculate G using output scores from temporal attention and link array

Require: L , A

```

1:  $idx \leftarrow 0$ 
2: for each  $t, s$  in  $L$  do
3:    $G \leftarrow []$ 
4:   if  $s > 0$  then
5:      $A' \leftarrow A[:, :, idx : idx + t]$ 
6:      $A' \leftarrow Expand(A', (n_h, N, t, s))$  {where  $Expand(tensor, dimension)$  expands tensor
       according to a given dimension}
7:      $A' \leftarrow Reshape(A', (n_h, N, ts))$ 
8:      $G \leftarrow G \cup A'$  {Append the respective temporal attention scores to  $G$ }
9:   end if
10:   $idx \leftarrow idx + t$ 
11: end for
12:  $G \leftarrow Stack(G)$  {Stack the list of tensors to construct tensor  $G$  of dimension  $(n_h, N, P)$ }

```

Given³ Q : the matrix of queries, K : keys of dimension d_k and V : values of dimension d_v , the scaled dot-product attention for the spatial layer is modified as:

$$Attention(Q, K, V, G) = \left(softmax\left(\frac{QK^T}{\sqrt{d_v}}\right) \odot G \right) V \quad (1)$$

In order to use the same CNP for multiple tasks with varying output vocabularies, we use multiple output embedding layers $OutputEmbedLayer_1, \dots, OutputEmbedLayer_{\tau_{len}}$, to generate the output embeddings for each task. At the final layer, we use multiple $(FC + Softmax)_1, \dots, (FC + Softmax)_{\tau_{len}}$ classification layers for each task. We also calculate a task embedding vector using τ and always start decoding using the task embedding vector.

3.3 Multi-task learning

In order to train a single a model simultaneously on mutiple tasks we used the HogWild training approach as described in [29]. Similar to the approach described in [30], the main process holds a

³For brevity, we reuse the notations of [1] in this description

global copy of the model. We create separate worker processes for each task, where each process maintains a local copy of a model. At each training iteration, each process starts by synchronizing its local model with the global copy. This is done through forward and backward propagation on its local copy and then copying the locally computed gradients to the global model asynchronously. Each process then calls the global model optimizer asynchronously to update the weights of the global model. Instead of storing the model in CPU as in [30] we always store the local copies across multiple GPUs.

4 Tasks and Setup

We evaluate the effectiveness of our proposed framework for four tasks of different domains (text, images, video): Image Captioning, Part-of-Speech (POS) tagging, Visual Question Answering (VQA) and Video-activity Recognition. The hyperparameter values used for n_h , d_{model} , N_Layers , d_k , d_v are same as that specified in Transformer base model [1]. For training, we always use cross-entropy loss with Adam optimizer [31] and schedule the learning rate using Noam scheduler [32] similar to [1]. In this section, we provide details on datasets used for the experimental analysis. We further elaborate on the model setup for each of these tasks.

Part-of-speech Tagging: To illustrate the task with only temporal modality ($t > 1$ & $s = 1$, where t and s are the temporal and spatial dimensions of the input respectively), we consider the part-of-speech (POS) tagging problem. Given an input sequence of words, the model should produce a sequence of POS tags corresponding to each word. We use Penn Tree-bank⁴ [33] which contains gold annotations on English WSJ articles by experts. We extract part-of-speech tag labels from the treebank and use splits 0-18 as training, 19-21 as development and 22-24 as test sets. During the encoding stage, each input sentence is passed to the language peripheral to generate a spatio-temporal tensor $x \in \mathbb{R}^{t \times 1 \times d_{model}}$, where t is the sequence length of subwords in the input. The $CNP_encode()$ function is then used to encode x into the temporal cache. Note that spatial cache is empty for text inputs, since $s = 1$. Thus the decoding stage, is same as that for Transformers to predict the sequence of POS tags.

Image Captioning: For a task with inputs containing only spatial modality ($t = 1$ & $s > 1$), we demonstrate the Image Captioning problem. The task requires a model to predict a text caption for a given image. We use the MSCOCO 2014 dataset [34] for training and present results on the COCO validation set. During the encoding stage, the input image is resized to dimension 224×224 and processed by the vision peripheral containing pre-trained ResNet-152 to produce image embeddings $x \in \mathbb{R}^{1 \times 49 \times d_{model}}$. x is then input to the $encode()$ function which populates corresponding spatial and temporal cache. The decoding stage uses $decode()$ function with output vocabulary size 25000 to generate the captions.

Visual Question Answering: For the task with inputs from multiple domain, such that each contains either spatial or temporal modality (either $t > 1$ & $s = 1$, or $t = 1$ & $s > 1$ for each input), we choose the task of visual question answering. This task has a question over an image as inputs and predicts the correct answer label. We use the recently introduced VQA v2.0 dataset [35] for this purpose and perform evaluation on the VQA test-dev set. This dataset was used in VQA 2018 challenge and contains over 1.1M question with 11.1M answers relating to images in MSCOCO dataset. All the images are resized to dimension 224×224 before training. The encoding stage of this task utilizes two peripherals: the vision peripheral is used to generate a tensor $x_1 \in \mathbb{R}^{1 \times 49 \times d_{model}}$ for the input image. The language peripheral is used to encode the questions into $x_2 \in \mathbb{R}^{t \times 1 \times d_{model}}$, where t is equal to the length of the subwords in the question. The $encode()$ function is called two times, first with x_1 and second with x_2 as input. Finally, the $decode()$ with output vocabulary size 3500 is to generate the answers as softmax probabilities in a single decoding step.

Video Activity Recognition: For tasks which contain both spatial and temporal modality in a single input ($t > 1$ & $s > 1$), we consider the action recognition task on videos. For this purpose, we use the HMDB dataset [36]. The dataset consists of over 5000 short length clips of real life actions with 51 classes. We present our results on train-test split 1. We use 16 frames per video and resize each of them to 224×224 . During the encoding stage, each frame of the video is passed through the vision peripheral to cumulatively generate a video encoding $x \in \mathbb{R}^{16 \times 49 \times d_{model}}$ which is then used as input

⁴<https://catalog.ldc.upenn.edu/LDC99T42>

	POS	Captioning	VQA			HMDB		#PARAMS
	Acc.	BLEU-1	Overall	Y/N	Num.	Other	Acc.	
State-of-the-art	97.44	77.2	63.2	80.3	42.8	55.8	59.4	-
Independent	95.61	70.04	55.31	74.09	35.17	46.35	55.29	450×10^6
Multi-task	95.44	69.04	55.76	75.49	35.64	46.08	54.44	149×10^6

Table 1: Performance of *OmniNet* on diverse set of tasks when trained individually as well as in multi-tasking setup.

to the *encode()* function. Finally, the *decode()* with output vocabulary size 51 is to predict the action as softmax probabilities in a single decoding step.

5 Results and Discussion

We present the evaluation on (a) Tasks of various modalities illustrated in 4 (b) Multi-tasking setup for these tasks (Table 1) (c) Reuse of the multi-task model for an unseen task (Figure 3). In addition, we also provide some ablation studies on the architecture (Table 2)

Performance of proposed architecture on individual tasks: We choose a set of four tasks with diverse input modalities and combinations as described in previous Section 4. We train the *OmniNet* model independently across each of the above tasks. Each of the tasks demonstrates unique capabilities of this generic architecture. More specifically, we compare our results with the following state-of-the-art⁵: POS tagging: [37]; image captioning & VQA: [16] and HMDB: [38]. The results of comparison with the state-of-art is shown in Table 1. While we do not obtain state-of-art performance in all the tasks, most of the tasks still attain comparable performance without any hyper-parameter optimization⁶. We do not focus on optimizing the performance for each of these specific tasks because our primary objective is to verify whether the model can achieve reasonable performance across all tasks using a single unified architecture. Although, we do not obtain the state-of-the-art performance in image captioning, VQA and HMDB, the fact that a single model can be used for tasks with varying combination of input modalities verifies the wide applicability of the proposed architecture. We further believe that, with more computational power, fine tuning the hyperparameters to these tasks should certainly result in improved performance. Its interesting to note that the model can be easily extended to new modalities like speech without any modification to the CNP as long as the inputs are provided in a spatio-temporal form.

Effect of training a diverse set of tasks together: We trained a single joint model on the four tasks mentioned above. Table 1 shows that the multi-task model attains similar performance as that of independent training across various tasks, while resulting in three times reduction in model size. That is, when a separate model is used for each task, we have a total of over 450×10^6 parameters. Whereas during multi-tasking, the model is much compressed with about 149×10^6 parameters, while achieving similar performance. Interestingly, the model is able to attend on spatio-temporal components of the inputs from different tasks and concurrently generate predictions across them, thus demonstrating the generalization capability of our architecture.

Impact of individual architectural components: In order to support different input domains, our architecture introduces spatial cache and link array components to the original *Transformer* architecture (which only consists of mechanisms to handle temporal data). We conducted an ablation study on each of these components to verify their importance across various tasks as shown in Table 2. The second row ablates the link array from our architecture i.e. removing the multiplication of G in Equation 1. The link array was designed to assist in tasks such as video, containing both spatial as well as temporal modality in a single input. The total number of spatial components becomes very large as number of frames in the video increases, thereby making it difficult to attend on various spatial regions throughout the video. Using link array the spatial attention layer can attend more on specific important frames in the video. Therefore, removal of link array leads to a huge reduction in performance in HMDB compared to other tasks, because they do not have both spatio-temporal modalities for any single input. Removal of spatial cache, on the other hand, has

⁵Since most of these tasks are popular challenges, we compare with state-of-the-art which are generically applicable for the respective task instead of the challenge dataset.

⁶We used the exact same hyper-parameters as specified in [1]

	POS	Captioning	VQA	HMDB
	Acc.	BLEU-1	Overall	Acc.
<i>OmniNet</i> Architecture	95.61	70.04	55.31	55.29
w/o link array	95.61	70.04	54.05	45.94
w/o spatial cache	95.61	0.0	39.86	10.91

Table 2: Ablation study on the effect of new architectural components (spatial cache and link array) introduced in the proposed architecture

significant effect on performance across all tasks containing spatial modality. Since, image captioning contains only spatial modality and hence the BLEU drops to 0 after ablation. On the other hand, VQA leverages both, spatial information from image and temporal information from question, retains some performance from the use of temporal cache. HMDB also contains spatially averaged frames in the temporal cache and hence retains some performance. Note that, POS tagging task is not affected by ablation of any of the components since it only has temporal modality in the input.

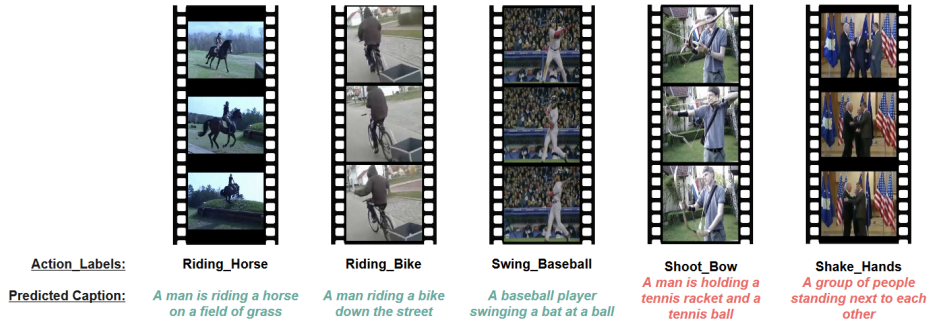


Figure 3: Results of zero-shot video captioning on a model pre-trained on image captioning and video recognition.

Towards zero-shot learning: reuse of pre-trained network for unseen task Sharing representations across multiple tasks provides the benefit to transfer of useful knowledge across multiple domains. Since, image and video are processed by the same vision peripheral, we conducted an experiment to see whether our model pre-trained on image captioning and video recognition can perform video captioning without any training on video captioning data (zero-shot learning). The results of the evaluation on randomly picked instances from the HMDB test split 1 are shown in Figure 3. Interestingly, the model performs quite well on related actions that were present in the COCO training; such as captions related to horse riding and baseball. Without training on any video captioning instance, the model could use the pre-trained information from image captioning to generate meaningful predictions, hence demonstrating the capability of the model to transfer knowledge across related multi-modal tasks. However, on concepts that are not present in the MSCOCO dataset, the model either describes the environment in the video or replaces with alternate known concepts. This case study, although not comprehensive, shows the capability of the model to learn shared representations and ability to transfer knowledge across domains. We believe, that adding more tasks and domains will lead to more interesting zero-shot learning results in future across a wide range of problems.

6 Conclusions and Future Work

We present a unified neural network architecture *OmniNet* capable of learning tasks with multiple inputs of varying modalities. Further, the architecture can be adopted for multi-task learning across any set of tasks containing spatio-temporal data. Our architecture attains comparable performance to the state-of-art on over four diverse tasks. Training these wide array of tasks concurrently provides a single compressed model while retaining similar performance. We further demonstrate that this shared model can learn robust representations from various spatio-temporal inputs which are reusable for unseen tasks.

We believe that this proposed architecture has wide applicability to any task with spatio-temporal inputs. To extend its usability, we would like to introduce new peripherals supporting more domains such as speech. We are keen on exploring other aspects to the data beyond temporal and spatial dimensions such as graphs and relational data.

References

- [1] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems 30*, pages 5998–6008. Curran Associates, Inc., 2017.
- [2] Yonghui Wu, Mike Schuster, Zhifeng Chen, Quoc V. Le, Mohammad Norouzi, Wolfgang Macherey, Maxim Krikun, Yuan Cao, Qin Gao, Klaus Macherey, Jeff Klingner, Apurva Shah, Melvin Johnson, Xiaobing Liu, Łukasz Kaiser, Stephan Gouws, Yoshikiyo Kato, Taku Kudo, Hideto Kazawa, Keith Stevens, George Kurian, Nishant Patil, Wei Wang, Cliff Young, Jason Smith, Jason Riesa, Alex Rudnick, Oriol Vinyals, Greg Corrado, Macduff Hughes, and Jeffrey Dean. Google’s neural machine translation system: Bridging the gap between human and machine translation. *CoRR*, abs/1609.08144, 2016.
- [3] Christian Szegedy, Sergey Ioffe, and Vincent Vanhoucke. Inception-v4, inception-resnet and the impact of residual connections on learning. *CoRR*, abs/1602.07261, 2016.
- [4] Eric Battenberg, Jitong Chen, Rewon Child, Adam Coates, Yashesh Gaur, Yi Li, Hairong Liu, Sanjeev Satheesh, David Seetapun, Anuroop Sriram, and Zhenyao Zhu. Exploring neural transducers for end-to-end speech recognition. *CoRR*, abs/1707.07413, 2017.
- [5] Rui Shu, Hung H. Bui, Hirokazu Narui, and Stefano Ermon. A dirt-t approach to unsupervised domain adaptation. *CoRR*, abs/1802.08735, 2018.
- [6] Lanqing Hu, Meina Kan, Shiguang Shan, and Xilin Chen. Duplex generative adversarial network for unsupervised domain adaptation. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018.
- [7] Jiquan Ngiam, Aditya Khosla, Mingyu Kim, Juhan Nam, Honglak Lee, and Andrew Y. Ng. Multimodal deep learning. In *Proceedings of the 28th International Conference on International Conference on Machine Learning*, ICML’11, pages 689–696, USA, 2011. Omnipress.
- [8] Ronan Collobert and Jason Weston. A unified architecture for natural language processing: Deep neural networks with multitask learning. In *Proceedings of the 25th International Conference on Machine Learning*, ICML ’08, pages 160–167, New York, NY, USA, 2008. ACM.
- [9] Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. Automatic differentiation in PyTorch. In *NIPS Autodiff Workshop*, 2017.
- [10] Melvin Johnson, Mike Schuster, Quoc V. Le, Maxim Krikun, Yonghui Wu, Zhifeng Chen, Nikhil Thorat, Fernanda Viégas, Martin Wattenberg, Greg Corrado, Macduff Hughes, and Jeffrey Dean. Google’s multilingual neural machine translation system: Enabling zero-shot translation. *Transactions of the Association for Computational Linguistics*, 5:339–351, 2017.
- [11] Daxiang Dong, Hua Wu, Wei He, Dianhai Yu, and Haifeng Wang. Multi-task learning for multiple language translation. In *ACL*, 2015.
- [12] M. L. Seltzer and J. Droppo. Multi-task learning in deep neural networks for improved phoneme recognition. In *2013 IEEE International Conference on Acoustics, Speech and Signal Processing*, pages 6965–6969, May 2013.
- [13] Kalpesh Krishna, Shubham Toshniwal, and Karen Livescu. Hierarchical multitask learning for ctc-based speech recognition. *CoRR*, abs/1807.06234, 2018.
- [14] Zhanpeng Zhang, Ping Luo, Chen Change Loy, and Xiaoou Tang. Facial landmark detection by deep multi-task learning. In *In ECCV. 94–108*, 2014.

- [29] Benjamin Recht, Christopher Re, Stephen Wright, and Feng Niu. Hogwild: A lock-free approach to parallelizing stochastic gradient descent. In J. Shawe-Taylor, R. S. Zemel, P. L. Bartlett, F. Pereira, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems 24*, pages 693–701. Curran Associates, Inc., 2011.
- [30] Volodymyr Mnih, Adrià Puigdomènech Badia, Mehdi Mirza, Alex Graves, Tim Harley, Timothy P. Lillicrap, David Silver, and Koray Kavukcuoglu. Asynchronous methods for deep reinforcement learning. In *Proceedings of the 33rd International Conference on International Conference on Machine Learning - Volume 48*, ICML’16, pages 1928–1937. JMLR.org, 2016.
- [31] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*, 2015.
- [32] Noam Shazeer and Mitchell Stern. Adafactor: Adaptive learning rates with sublinear memory cost. *CoRR*, abs/1804.04235, 2018.
- [33] Mitchell Marcus, Grace Kim, Mary Ann Marcinkiewicz, Robert MacIntyre, Ann Bies, Mark Ferguson, Karen Katz, and Britta Schasberger. The penn treebank: Annotating predicate argument structure. In *Proceedings of the Workshop on Human Language Technology, HLT ’94*, pages 114–119, Stroudsburg, PA, USA, 1994. Association for Computational Linguistics.
- [34] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C. Lawrence Zitnick. Microsoft coco: Common objects in context. In David Fleet, Tomas Pajdla, Bernt Schiele, and Tinne Tuytelaars, editors, *Computer Vision – ECCV 2014*, pages 740–755, Cham, 2014. Springer International Publishing.
- [35] Yash Goyal, Tejas Khot, Douglas Summers-Stay, Dhruv Batra, and Devi Parikh. Making the V in VQA matter: Elevating the role of image understanding in Visual Question Answering. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017.
- [36] H. Kuehne, H. Jhuang, E. Garrote, T. Poggio, and T. Serre. HMDB: a large video database for human motion recognition. In *Proceedings of the International Conference on Computer Vision (ICCV)*, 2011.
- [37] Drahomíra “johanka” Spoustová, Jan Hajič, Jan Raab, and Miroslav Spousta. Semi-supervised training for the averaged perceptron POS tagger. In *Proceedings of the 12th Conference of the European Chapter of the ACL (EACL 2009)*, pages 763–771, Athens, Greece, March 2009. Association for Computational Linguistics.
- [38] Karen Simonyan and Andrew Zisserman. Two-stream convolutional networks for action recognition in videos. In Z. Ghahramani, M. Welling, C. Cortes, N. D. Lawrence, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems 27*, pages 568–576. Curran Associates, Inc., 2014.