

A Priority Map for Vision-and-Language Navigation with Trajectory Plans and Feature-Location Cues

JASON ARMITAGE
University of Zurich
Switzerland

LEONARDO IMPETT
University of Cambridge
UK

RICO SENNRICH
University of Zurich
Switzerland

ABSTRACT

In a busy city street, a pedestrian surrounded by distractions can pick out a single sign if it is relevant to their route. Artificial agents in outdoor Vision-and-Language Navigation (VLN) are also confronted with detecting supervisory signal on environment features and location in inputs. To boost the prominence of relevant features in transformer-based systems without costly preprocessing and pretraining, we take inspiration from priority maps - a mechanism described in neuropsychological studies. We implement a novel priority map module and pretrain on auxiliary tasks using low-sample datasets with high-level representations of routes and environment-related references to urban features. A hierarchical process of trajectory planning - with subsequent parameterised visual boost filtering on visual inputs and prediction of corresponding textual spans - addresses the core challenge of cross-modal alignment and feature-level localisation. The priority map module is integrated into a feature-location framework that doubles the task completion rates of standalone transformers and attains state-of-the-art performance for transformer-based systems on the Touchdown benchmark for VLN. Code and data are referenced in Appendix C.

1 INTRODUCTION

Navigation in the world depends on attending to relevant cues at the right time. A road user in an urban environment is presented with billboards, moving traffic, and other people - but at an intersection will pinpoint a single light to check if it contains the colour red (Gottlieb et al., 2020; Shinoda et al., 2001). An artificial agent navigating a virtual environment of an outdoor location is also presented with a stream of linguistic and visual cues. Action selections that move the agent closer to a final destination depend on the prioritisation of references that are relevant to the point in the trajectory. In the first example, human attention is guided to specific objects by visibility and the present objective of crossing the road. At a neurophysiological level, this process is mediated by a priority map - a neural mechanism that guides attention by matching low-level signals on salient objects with high-level signals on task goals. Prioritisation in humans is enhanced by combining multimodal signals and integration between linguistic and visual information (Ptak, 2012; Cavicchio et al., 2014). The ability to prioritise improves as experience of situations and knowledge of environments increases (Zelinsky and Bisley, 2015; Tatler et al., 2011).

We introduce a priority map module for Vision-and-Language Navigation (PM-VLN) that is pretrained to guide a transformer-based architecture to prioritise relevant information for action selections in navigation. In contrast to pretraining on large-scale datasets with generic image-text pairs (Su et al., 2020), the PM-VLN module learns from small sets of samples representing trajectory plans and urban features. Our proposal is founded on observation of concentrations in location deictic terms and references to objects with high visual salience in inputs for VLN. Prominent features in the environment pervade human-generated language navigation instructions. Road network types (“intersection”), architectural features (“awning”), and transportation (“cars”) all appear with high frequency in linguistic descriptions of the visual appearance of urban locations. Learning to combine information in the two modalities relies on synchronising temporal sequences of varying lengths. We utilise references to entities as a signal for a process of cross-modal prioritisation that addresses this requirement.

Our module learns over both modalities to prioritise timely information and assist both generic vision-and-language and custom VLN transformer-based architectures to complete routes (Li et al., 2019; Zhu et al., 2021). Transformers have contributed to recent proposals to conduct VLN, Visual Question Answering, and other multimodal tasks - but are associated with three challenges: 1) Standard architectures lack mechanisms that address the challenge of temporal synchronisation over linguistic and visual inputs. Pretrained transformers perform well in tasks on image-text pairs but are challenged when learning over sequences without explicit alignments between modalities (Lin and Wang, 2020). 2) Performance is dependent on pretraining with large sets of image-text pairs and a consequent requirement for access to enterprise-scale computational resources (Majumdar et al., 2020; Suglia et al., 2021). 3) Visual learning relies on external models and pipelines - notably for object detection (Li et al., 2020; Le et al., 2022). The efficacy of object detection for VLN is low in cases where training data only refer to a small subset of object types observed in navigation environments.

We address these challenges with a hierarchical process of trajectory planning with feature-level localisation and low-sample pretraining on in-domain data. We use discriminative training on two auxiliary tasks that adapt parameters of the PM-VLN for the specific challenges presented by navigating routes in outdoor environments. High-level planning for routes is enabled

by pretraining for trajectory estimation on simple path traces ahead of a second task comprising multi-objective cross-modal matching and location estimation on urban landmarks. Data in the final evaluation task represent locations and trajectories in large US cities and present an option to leverage real-world resources in pretraining. Our approach builds on this opportunity by sourcing text, images, coordinates, and path traces from the open web and the Google Directions API where additional samples may be secured at low cost in comparison to human generation of instructions.

This research presents four contributions to enhance transformer-based systems on outdoor VLN tasks:

- **Priority map module** Our novel PM-VLN module conducts a hierarchical process of high-level alignment of textual spans with visual perspectives and feature-level operations to enhance and localise inputs during navigation (see Figure 3).
- **Trajectory planning** We propose a new method for aligning temporal sequences in VLN comprising trajectory estimation on path traces and subsequent predictions for the distribution of linguistic descriptions over routes.
- **Two in-domain datasets and training strategy** We introduce a set of path traces for routes in two urban locations (TR-NY-PIT-central) and a dataset consisting of textual summaries, images, and World Geodetic System (WGS) coordinates for landmarks in 10 US cities (MC-10). These resources enable discriminative training of specific components of the PM-VLN on trajectory estimation and multi-objective loss for a new task that pairs location estimation with cross-modal sentence prediction.
- **Feature-location framework** We design and build a framework (see Figure 2) to combine the outputs from the PM-VLN module and cross-modal embeddings from a transformer-based encoder. The framework incorporates components for performing self-attention, combining embeddings, and predicting actions with maxout activation.

2 BACKGROUND

In this section we define the Touchdown task and highlight a preceding challenge of aligning and localising over linguistic and visual inputs addressed in our research. A summary of the notation used below and in subsequent sections is presented in Appendix A.

Touchdown Navigation in the Touchdown benchmark ϕ_{VLN} is measured as the completion of N predefined trajectories by an agent in an environment representing an area of central Manhattan. The environment is represented as an undirected graph composed of nodes O located at WGS latitude/longitude points. At each step t of the sequence $\{1, \dots, T\}$ that constitute a trajectory, the agent selects an edge ξ_t to a corresponding node. The agent’s selection is based on linguistic and visual inputs. A textual instruction τ composed of a varying number of tokens describes the overall trajectory. We use ς to denote a span of tokens from τ that corresponds to the agent’s location in the trajectory. Depending

on the approach, ς can be the complete instruction or a selected sequence. The visual representation of a node in the environment is a panorama drawn from a sequence $Route$ of undetermined length. The agent receives a specific perspective ψ of a panorama determined by the heading angle \angle between (o_1, o_2) . Success in completing a route is defined as predicting a path that ends at the node designated as the goal - or one directly adjacent to it.

In a supervised learning paradigm (see a) in Figure 1), an embedding e_η is learned from inputs ς_t and ψ_t . The agent’s next action is a classification over e_η where the action α_t is one of a class drawn from the set $A\{Forward, Left, Right, Stop\}$. Predictions $\alpha_t = Forward$ and $\alpha_t = \{Left, Right\}$ result respectively in a congruent or a new \angle at edge ξ_{t+1} . A route in progress is terminated by a prediction $\alpha_t = Stop$.

Align and Localise We highlight in Figure 1 a preceding challenge in learning cross-modal embeddings. As in real-world navigation, an agent is required to align and match cues in instructions with its surrounding environment. A strategy in human navigation is to use entities or landmarks to perform this alignment Cavicchio et al. (2014). In the Touchdown benchmark, a relationship between sequences τ and $Route$ is assumed from the task generation process outlined in Chen et al. (2019) - but the precise alignment is not known. We define the challenge as one of aligning temporal sequences $\tau = \{\varsigma_1, \varsigma_2, \dots, \varsigma_n\}$ and $Route = \{\psi_1, \psi_2, \dots, \psi_n\}$ with the aim of generating a set of cross-modal embeddings E_η where referenced entities correspond. At a high level, this challenge can be addressed by an algorithm q that maximises the probability P of detecting S signal in a set of inputs. This algorithm is defined as

$$q(\theta)((X_t)) = q(\theta)(X_t) \rightarrow \max \left[\int_{\mathcal{X}} p(X_t|\theta)s(X_t) \right] \quad (1)$$

where θ is a parameter $\theta \in \Theta_1$ and \mathcal{X} is the data space.

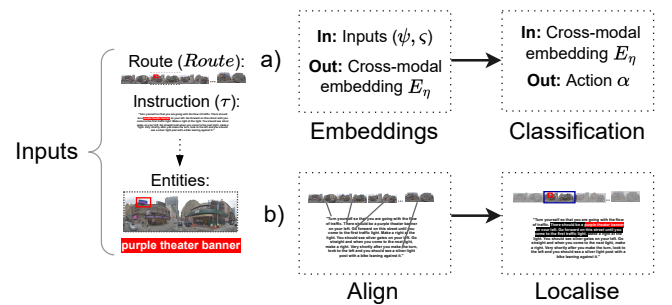


Figure 1: Outline of VLN as a supervised classification task a). Linguistic and visual inputs both refer to entities indicated in red. We address a challenge to align and localise over unsynchronised inputs b) by focusing on entities represented in both modalities.

In the Touchdown benchmark, linguistic and visual inputs are of the form $0 \leq |\tau| \leq n$ and $0 \leq |Route| \leq n$ where $len(\tau) \neq len(Route)$. The task then is to maximise the probability of detecting signal in the form of corresponding entities over the sequences τ and $Route$, which in turn is the product of

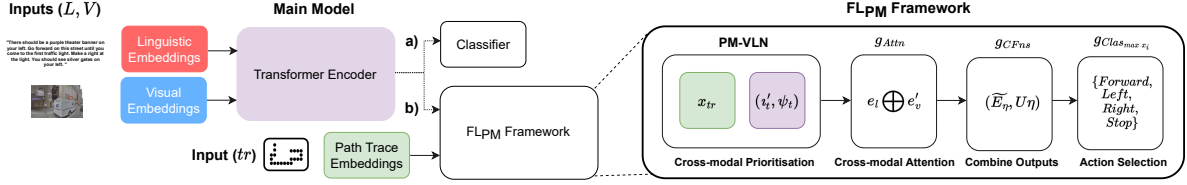


Figure 2: Prior work on transformer-based systems for VLN follows the above pipeline from inputs to the main model concluding with a) a classifier to predict actions. We propose a feature-location framework (FL_{PM}) to enhance the performance of a main model as in b). Here path traces are an additional input to assist the PM-VLN to align linguistic and visual sequences. Submodule g_{CFNs} combines embeddings from the main model U_η and the PM-VLN \bar{E}_η ahead of action prediction with maxout activation.

probabilities over pairings ζ_t and ψ_t presented at each step:

$$g(X_t) \rightarrow \max_{\text{subject to}} P[\tau, Route] = \prod p_{x_\zeta x_\psi} \quad (2)$$

3 METHOD

We address the challenge of aligning and localising over sequences with a computational implementation of cross-modal prioritisation. Diagnostics on VLN systems have placed in question the ability of agents to perform cross-modal alignment (Zhu et al., 2022). Transformers underperform in problems with temporal inputs where supervision on image-text alignments is lacking (Chen et al., 2020). This is demonstrated in the case of Touchdown where transformer-based systems complete less than a quarter of routes. Our own observations of lower performance when increasing the depth of transformer architectures motivates moving beyond stacking blocks to an approach that compliments self-attention.

Our PM-VLN module modulates transformer-based encoder embeddings in the main task ϕ_{VLN} using a hierarchical process of operations and leveraging prior learning on auxiliary tasks (ϕ_1, ϕ_2) (see Figure 3). In order to prioritise relevant information, a training strategy for PM-VLN components is designed where training data contain samples that correspond to the urban grid type and environment features in the main task. The datasets required for pretraining contain less samples than other transformer-based VLN frameworks (Zhu et al., 2021; Majumdar et al., 2020) and target only specific layers of the PM-VLN module. The pretrained module is integrated in a novel feature-location framework FL_{PM} shown in Figure 2. Subsequent components in the FL_{PM} combine cross-modal embeddings from the PM-VLN and a main transformer model ahead of predicting an action.

3.1 Feature-location Framework with a Priority Map Module

Prior work on VLN agents has demonstrated reliance for navigation decisions on environment features and location-related references (Zhu et al., 2021). In the definition of ϕ_{VLN} above, we consider this information as the supervisory signal contained in both sets of inputs $(x_\zeta, x_\psi)_t$. As illustrated in Figure 2, our PM-VLN module is introduced into a framework FL_{PM}. This framework takes outputs from a transformer-based main model Enc_{Trans} together with path traces ahead of cross-modal pri-

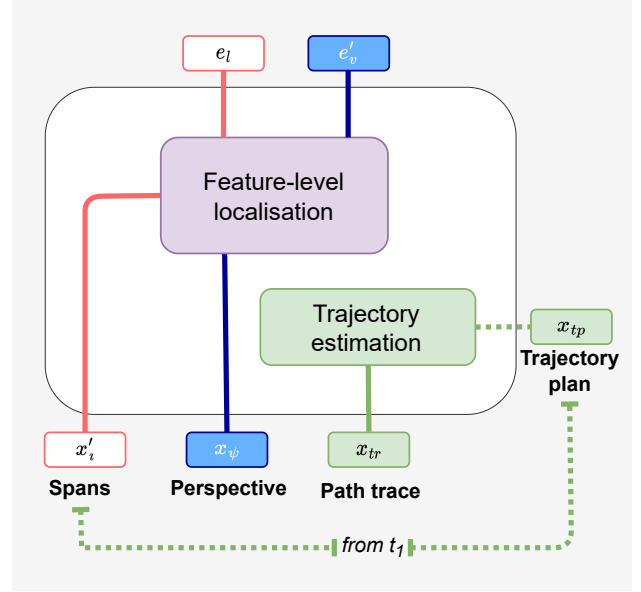


Figure 3: A Priority Map module performs a hierarchical process of high-level trajectory planning and feature-level localisation. Sub-modules inside the white box are learned together and a helper function generates a trajectory plan to predict spans from step t_1 .

oritisation and classification with maxout activation $Clas_{max} x_i$. Inputs for Enc_{Trans} comprise cross-modal embeddings \bar{e}_η proposed by Zhu et al. (2021) and a concatenation of perspectives up to the current step ψ_{cat}

$$Clas_{max} x_i [y_j | z'] = \bar{d}(PM-VLN(\{g(x_i), (tr_t, t'_i, \psi_t)\}_{i=1}^n) + Enc_{Trans}(\{g(x_i), (\bar{e}_\eta, \psi_{cat})\}_{i=1}^n)) \quad (3)$$

where tr_t is a path trace, z' is the concatenation of the outputs of the two encoders, and \bar{d} is a dropout operation.

3.1.1 Priority Map Module Priority maps are described in the neuropsychological literature as a mechanism that modulates sensory processing on cues from the environment. Saliency deriving from the physical aspects of objects in low-level processing is mediated by high-level signals for the relevance of cues to task goals (Fecteau and Munoz, 2006; Itti and Koch, 2000; Zelinsky and Bisley, 2015). Cortical regions that form the location of

these mechanisms are associated with the combined processing of feature- and location-based information (Bisley and Mirpour, 2019; Hayden and Gallant, 2009). Prioritisation of items in map tasks with language instructions indicate an integration between linguistic and visual information and subsequent increases in salience attributed to landmarks (Cavicchio et al., 2014).

Our priority map module (PM-VLN) uses a series of simple operations to approximate the prioritisation process observed in human navigation. These operations avoid dependence on initial tasks such as object detection. Alignment of linguistic and visual inputs is enabled by trajectory estimation on simple path traces forming high-level representations of routes and subsequent generation of trajectory plans. Localisation consists of parameterised visual boost filtering on the current environment perspective ψ_t and cross-modal alignment of this view with selected spans from subsequent alignment (see Algorithm 1). This hierarchical process compliments self-attention by accounting for the lack of a mechanism in transformers to learn over unaligned temporal sequences. A theoretical basis for cross-modal prioritisation is presented below.

Algorithm 1 Priority Map Module

Input: Datasets $\mathcal{D}_{\phi_1}, \mathcal{D}_{\phi_2}$, and $\mathcal{D}_{\phi_{VLN}}$ with inputs (x_l, x_v) for tasks Φ . Initial parameters in all layers at $\Theta_j^l \sim Normal(\mu_j, \sigma_j)$.

Output: (e_l, e'_v)

while not converged **do**

for x_{tr_i} in ϕ_1 **do**

$\Theta'_{gPMTp} \leftarrow g_{\phi_1}(X_i, \Theta)$.

end for

end while

while not converged **do**

for (x_{l_i}, x_{v_i}) in ϕ_2 **do**

$\Theta'_{gPMF} \leftarrow g_{\phi_2}(X_i, \Theta)$.

end for

end while

while not converged **do**

 Sample x_{tr_t} from D^{Train} .

$x_{tp_t} \leftarrow g_{PMTp}(x_{tr_t})$.

 Sample (t'_t, ψ_t) from D^{Train} .

$e_v \leftarrow g_{USM}(\psi_t)$.

$e'_v \leftarrow g_{VBF}(e_v)$.

$e_l \leftarrow g_{PrL}(g_{Cat}(t'_t, e'_v))$.

end while

return (e_l, e'_v)

High-level trajectory estimation Alignment over linguistic and visual sequences is formulated as a task of predicting a set of spans from the instruction that correspond to the current step. This process starts with a submodule g_{PMTp} that estimates a count cnt of steps from a high-level view on the route (see Figure 4). Path traces - denoted as tr_T - are visual representations of trajectories generated from the coordinates of nodes. At t_0 in tr_T initial

spans in the instruction are assumed to align with the first visual perspective. From step t_1 , a submodule containing a pretrained ConvNeXt Tiny model (Liu et al., 2022) updates an estimate of the step count in cnt_{tr_T} . A trajectory plan tp_t is a Gaussian distribution of spans in τ within the interval $[x_{left}, x_{right}]$. At each step, samples from this distribution serve as a prediction for relevant spans. The final output t'_t is the predicted span t_t combined with t_{t-1} .

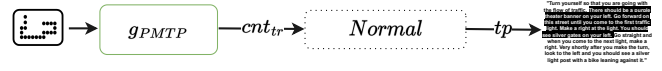


Figure 4: Submodule g_{PMTp} estimates a step count (cnt_{tr}) on a path trace. A trajectory plan (tp) is a Gaussian distribution ($Normal$) over the instruction and predicts a span for every step t_t . This is concatenated with the span predicted for the previous step.

Feature-level localisation Predicted spans are passed with ψ_t to a submodule g_{PMF} that is pretrained on cross-modal matching in ϕ_2 (see Figure 5). Feature-level operations commence with visual boost filtering. Let $Conv_{VBF}$ be a convolutional layer with a kernel κ and weights W that receives as input ψ_t . In the first operation g_{USM} , a Laplacian of Gaussian kernel κ_{LoG} is applied to ψ_t . The second operation g_{VBF} consists of subtracting the output e_v from the original tensor ψ_t :

$$g_{VBF}(e_v) = (\lambda - 1)(e_v) - g_{USM}(\psi_t) \quad (4)$$

where λ is a learned parameter for the degree of sharpening.

A combination of g_{USM} and g_{VBF} is equivalent to adaptive sharpening of details in an image with a Laplacian residual (Bayer, 1986). Here operations are applied directly to e_v and adjusted at each update of $W_{Conv_{VBF}}$ with a parameterised control $\beta\lambda$. In the simple and efficient implementation from Carranza-Rojas et al. (2019), σ in the distribution $LoG(\mu_j, \sigma_j)$ is fixed and the level of boosting is reduced to a single learned term

$$\Delta z(x_1, x_2) = \beta\lambda \left(\sum_j (AA'_{\kappa_{ij}} - A_{W_{\kappa_{ij}}}) z \right) \quad (5)$$

where A_W is a matrix of parameters and AA' is the identity.

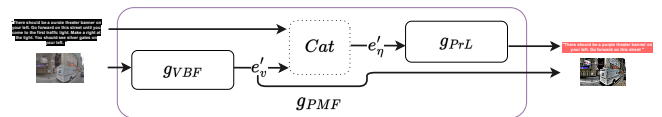


Figure 5: Submodule g_{PMF} commences feature-level operations by boosting visual features in the perspective. The next operation (Cat) is a concatenation of the output from g_{VBF} and the linguistic output t'_t from the alignment process above. A precise prediction for the relevant span e_l is returned by g_{PrL} .

Selection of a localised span e_l proceeds with a learned cross-modal embedding e'_η composed of e'_v and the linguistic output t'_t from the preceding alignment operation. A binary prediction over

this linguistic pair is performed on the output hidden state from a single-layer LSTM, which receives e'_η as its input sequence. Function g_{PrL} returns a precise localisation of relevant spans w.r.t. prominent features in the perspective:

$$g_{PrL}(e_l) = g_{Cat}(t'_l, e'_v) \triangleq \begin{cases} 0, & \text{if } \langle w, x \rangle + b < 0 \\ 1, & \text{otherwise} \end{cases} \quad (6)$$

Pretraining Strategy A data-efficient pretraining strategy for the PM-VLN module consists of pretraining submodules of the PM-VLN on auxiliary tasks (ϕ_1, ϕ_2). We denote the two datasets for these tasks as ($\mathcal{D}_{\phi_1}, \mathcal{D}_{\phi_2}$) and a training partition as \mathcal{D}^{Train} (see Appendix B for details). In ϕ_1 , the g_{PMTp} submodule is pretrained on TR-NY-PIT-central - a new set of path traces. Path traces in $\mathcal{D}_{\phi_1}^{Train}$ are generated from 17,000 routes in central Pittsburgh with a class label for the step count in the route. The distribution of step counts in $\mathcal{D}_{\phi_1}^{Train}$ is 50 samples for routes with ≤ 7 steps and 300 samples for routes with > 7 steps (see Appendix B). During training, samples from $\mathcal{D}_{\phi_1}^{Train}$ are presented in standard orientation for 20 epochs and rotated 180° ahead of a second round of training. This rotation policy is preferred following empirical evaluation using standalone versions of the g_{PMTp} submodule receiving two alternate preparations of $\mathcal{D}_{\phi_1}^{Train}$ with random and 180° rotations. Training is formulated as multiclass classification with cross-entropy loss on a set of $M=66$ classes

$$g_{\phi_1}(x_{tr}, \Theta) = B_0 + \underset{i}{\operatorname{argmax}} \sum_{j=1}^M B_i(x_{tr}, W_j) \quad (7)$$

where a class is the step count, B is the bias, and i is the sample in the dataset.

Pretraining on ϕ_2 for the feature-level localisation submodule g_{PMF} is conducted with the component integrated in the framework FL_{PM} and the new MC-10 dataset. Samples in $\mathcal{D}_{\phi_2}^{Train}$ consist of 8,100 landmarks in 10 US cities. To demonstrate the utility of open source tools in designing systems for outdoor VLN, the generation process leverages free and accessible resources that enable targeted querying. Entity IDs for landmarks sourced from the Wikidata Knowledge Graph are the basis for downloading textual summaries and images from the MediaWiki and Wikimedia APIs. Additional details on MC-10 are available in Appendix B. The aim in generating the MC-10 dataset is to optimise $\Theta_{g_{PMF}}$ such that features relating to $Y_{\phi_{VLN}}$ are detected in inputs $X_{\phi_{VLN}}$. We opt for open A multi-objective loss for ϕ_2 consists of cross-modal matching over the paired samples (x_l, x_v) - and a second objective comprising a prediction on the geolocation of the entity. In the first objective, g_{PMF} conducts a binary classification between the true x_l matching x_v and a second textual input selected at random from entities in the mini-batch. A limit of 540 tokens is set for all textual inputs and the classification in g_{PMF} is performed on the first sentence for each entity. Parameters $\Theta_{g_{PMF}}$ are saved and used subsequently for feature-level localisation in ϕ_{VLN} .

3.1.2 Cross-modal Attention and Action Prediction on Combined Outputs Resuming operations subsequent to the PM-VLN, outputs e'_{v_t} from $Conv_{VBF}$ are passed together with e_{l_t} to a VisualBERT embedding layer. Embeddings for both modalities are then processed by 4 transformer encoder layers with a hidden size of 256 and self-attention \oplus is applied to learn alignments between the pairs

$$\tilde{e}_\eta = \oplus(e_l \iff e'_v) = \operatorname{Soft} \left(\sum_{k=1}^{\mathcal{E}} \mathcal{M}_k \mathbb{L}(\mathcal{E}_k, \tilde{\mathcal{E}}_k) \right) \quad (8)$$

where Soft is the softmax function, k is the number of elements in the inputs, $\mathcal{M}_{k=1}$ is a masked element over the cross-modal inputs, \mathbb{L} is the loss, \mathcal{E}_k is an element in the input modality, and $\tilde{\mathcal{E}}_k$ is the predicted element. Cross-modal embeddings resulting from this attention operation are processed by concatenating over layer outputs $g(\tilde{e}_\eta) = (\tilde{e}_\eta^1, \tilde{e}_\eta^2, \tilde{e}_\eta^3, \tilde{e}_\eta^4)$.

Architectural and embedding selections for our frameworks aim to enable comparison with benchmark systems on ϕ_{VLN} . The Enc_{Tran} in the best performing framework uses a standard VisualBERT encoder with a hidden size of 256 and 4 layers and attention heads. As noted above, inputs for Enc_{Tran} align with those used in prior work (Zhu et al., 2021).

A submodule g_{CFns} combines U_η from \mathcal{L}^4 of the Enc_{Tran} and outputs from the cross-modal attention operation $g(\tilde{E}_\eta)$ ahead of applying dropout. Predictions for navigation actions are the outputs of a classifier block consisting of linear layers with maxout activation. Maxout activation in a block composed of linear operations takes the $\max z_{ij}$ where z_{ij} are the product of $x_{ij} W_{n^*}$ for k layers. In contrast to ReLU, the activation function is learned and prevents unit saturation associated with performing dropout (Goodfellow et al., 2013). We compare a standard classifier to one with $\max x_i$ in Table 2. Improvements with $\max x_i$ are consistent with a requirement to offset variance when training with the high number of layers in the full FL_{PM} framework.

3.2 Theoretical Basis

This section provides a theoretical basis for a hierarchical process of cross-modal prioritisation that optimises attention over linguistic and visual inputs. In this section we use q to denote this process for convenience. During the main task ϕ_{VLN} , q aligns elements in temporal sequences τ and $Route$ and localises spans and visual features w.r.t. a subset of all entities Ent in the routes:

$$q_{PM} = \|x_l - x_v\| \underset{\text{subject to}}{\rightarrow} \max P_{D_{Ent}}[\tau, Route] \leq R \quad (9)$$

Inputs in ϕ_{VLN} consist of a linguistic sequence τ and a visual sequence $Route$ or each trajectory j in a set of trajectories. As a result of the process followed by Chen et al. (2019) to create the Touchdown task, these inputs conform to the following definition.

Definition 1 (Sequences refer to corresponding entities). *At each step in j , $|x_l|$ and $|x_v|$ are finite subsequences drawn from τ_j and $Route_j$ that refer to corresponding entities appearing in the trajectory $ent_j \subset Ent$.*

In order to simplify the notation, these subsequences are denoted in this section as x_l and x_v . Touchdown differs from other outdoor navigation tasks (Hermann et al., 2020) in excluding supervision on the alignment over cross-modal sequences. Furthermore $len(\tau_j) \neq len(Route_j)$ and there are varying counts of subsequences and entities in trajectories. In an approach to ϕ_{VLN} formulated as supervised classification, an agent’s action at each step $\alpha_t \equiv$ classification $c_t \in \{0, 1\}$ where c is based on corresponding ent_t in the pair $(x_l, x_v)_t$. The likelihood that c_t is the correct action depends in turn on detecting S signal in the form of ent_t from noise in the inputs. The objective of q then is to maximise P_S for each point in the data space.

The process defining q is composed of multiple operations to perform two functions of high-level alignment g_{Align} and localisation g_{Loc} . At the current stage stg , function g_{Align} selects one set of spans $\varphi_{stg} \in (\varphi_1, \varphi_2, \dots, \varphi_n)$

$$\text{where } stg \begin{cases} \text{Start, if } t = 0 \\ \text{End, if } t = -1 \\ \forall stg_{other}, n \in N \in \sum_{n=1}^{n_1} > t_{-1} \text{ otherwise.} \end{cases}$$

This is followed by the function g_{Loc} , which predicts one of $\zeta_{scnt_0} \vee \zeta_{scnt_{0-1}}$ as the span ζ relevant to the current trajectory step $scnt$

$$\text{where } scnt \begin{cases} scnt_0, \text{ if } (\tau, \psi_t) = 0 \\ scnt_{0-1}, \text{ otherwise.} \end{cases}$$

We start by describing the learning process when the agent in ϕ_{VLN} is a transformer-based architecture $Enc + Clas$ excluding an equivalent to q (e.g. VisualBERT in Table 1 of the main report). $Enc + Clas$ is composed of two core subprocesses: cross-modal attention to generate representations $q(\bigoplus(L \iff \tilde{V}))$ and a subsequent classification $Clas(\tilde{e}_n')$.

Definition 2 (Objective in $Enc + Clas$). *The objective $Obj_1(\theta)$ for algorithm $q(\bigoplus(L \iff \tilde{V}))$, where L and V are each sequences of samples $\{x_1, x_2, \dots, x_n\}$, is the correspondence between samples x_l and x_v presented at step t in $\sum_{i=1}^n t_i = t_1 + t_2, \dots + t_n$.*

It is observed that in the learning process for $Enc + Clas$, any subprocesses to align and localise finite sequences x_l and x_v w.r.t. ent_j are performed as implicit elements in the process of optimising $Obj_1(\theta)$. In contrast the basis for the hierarchical learning process enabled by our framework FL_{PM} - which incorporates q_{PM} with explicit functions for these steps - is given in Theorem 1.

Theorem 1. *Assuming x_l and x_v conform to Definition 1 and that $\forall x \in L \exists x \in V$, an onto function $g_{Map} = mx + b, m \neq 0$ exists such that:*

$$g_{Map}(x_l, x_v) \rightarrow \max \left[ent_j^{(x_l, x_v)} \in Ent \right] \quad (10)$$

In this case, additional functions g_{Align} and g_{Loc} - when imple-

mented in order - maximise g_{Map} :

$$\max P_{D_{ent_j}} = \max g_{Map}(x_l, x_v) \xrightarrow{\text{subject to}} \overrightarrow{(g_{Align}, g_{Loc}, g_{Map})} \forall ent_j^{(x_l, x_v)} \in L_j \cap V_j \quad (11)$$

Remark 1 *Let $P(\max g_{Map})$ in Theorem 1 be the probability of optimising g_{Map} such that the number of pairs $N^{(x_l, x_v)}$ corresponding to $ent_j \in L_j \cap V_j$ is maximised. It is noted that $N^{(x_l, x_v)}$ is determined by all possible outcomes in the set of cases $\{(x_l, x_v) \iff ent_j, (x_l, x_v) \not\iff ent_j, x_l \not\iff x_v\}$. As the sequences of instances i in x_l, x_v and ent_j are forward-only, it is also noted that $N_{t+1}^{(x_l, x_v)} < N_t^{(x_l, x_v)}$ if $ent_i \notin x_{li}, ent_i \notin x_{vi}$, or $ent_i^{x_l} \neq ent_i^{x_v}$. By definition, $N_{t+1}^{(x_l, x_v)} > N_t^{(x_l, x_v)}$ if $P(ent_i = x_{li} = x_{vi})$ - where the latter probability is s.t. processes performed within finite computational time $CT(n)$ - which implies that $P(\max g_{Map}) | P(ent_i = x_{li} = x_{vi})$.*

Remark 2. *Following on from Remark 1, $CT(n^{P(ent_i = x_{li} = x_{vi})})$ when q contains g_t , and function $g_t(\max(N^{(x_l, x_v)} \iff ent_j \in L_j \cap V_j))$, where $g_t \in G < CT(n^{P(ent_i = x_{li} = x_{vi})})$ when q does not contain $g_t < CT(n^{P(ent_i = x_{li} = x_{vi})})$ when q contains g_t , and function $g_t(\max(N^{(x_l, x_v)} \not\iff ent_j \in L_j \cap V_j))$.*

Discussion In experiments, we expect from Remark 1 that results on ϕ_{VLN} for architectures such as $Enc + Clas$ - which exclude operations equivalent to those undertaken by the onto function g_{Map} - will be lower than the results for a framework FL_{PM} over a finite number of epochs. We observe this in Table 1 when comparing the performance of respective standalone and + FL_{PM} for VisualBERT and VLN Transformer systems. Poor results for variants (a) and (h) in Tables 2 and 3 in comparison to $FL_{PM} + VisualBERT(4l)$ also support the expectation set by Remark 2 that performance will be highly impacted in an architecture where operations in g_{Map} increase the number of misalignments.

Proof of Theorem 1 *We use below a^* for a generic transformer-based system that predicts α on (L, V) , ∇x for gradients, and Θ^{a^*} to denote $\Theta^{Enc+Clas} \vee \Theta^{Enc+q}$. Let sequence $x_l = [ent_1, ent_2, \dots, ent_{n_1}]$ and sequence $x_v = [ent_1, ent_2, \dots, ent_{n_2}]$, where n_1 and n_2 are unknown. We note that at any point during learning, $P_S(x_l, x_v)$ is spread unevenly over ent_j in relation to $\Theta^{a^*} \approx \mathcal{X}$.*

Propositions *We start with the case that $\exists ent_j : ent^{(x_l)}$ and $ent^{(x_v)}$. Here $CT(n^{Ent \in L \cap V})$ for $\Theta^{a^*+g_t} < CT(n^{Ent \in L \cap V})$ for Θ^{a^*} where g_t accounts for $\Delta(Len_1, Len_2)$. We next consider the case where $\nexists ent_j : ent^{(x_l)} \vee ent^{(x_v)}$. Where $\nexists g_{Loc}$ then $P_S^{(x_l, x_v)} < \exists g_{Loc} P_S^{(x_l, x_v)}$. We conclude with the case where $\exists Ent : x_l \vee x_v$. In $P_S^{A^*} ent^{(x_l)} \bigoplus ent^{(x_v)}$ when $ent^{(x_l)} \neq ent^{(x_v)}$.*

As $(Ent_L, Ent_V) \Rightarrow Ent, \Theta^{a^} \approx \max(N^{(x_l, x_v)}) \in \mathcal{X}$. $P_S^{(x_l, x_v)}$ where $ent_i = x_{li} = x_{vi} > ent_i \in \Theta^{a^*} \approx \max(N^{(x_l, x_v)})$. Furthermore $P \exists ent \in Ent \approx (ent_i) > \nexists ent \neq ent_i$. Therefore slope ∇x increases and $CT(n^{Ent \in L \cap V})$ for $\Theta^{a^*+q} < CT(n^{Ent \in L \cap V})$.*

		Development			Test		
		TC↑	SPD↓	SED↑	TC↑	SPD↓	SED↑
Inputs (L, V) (non-transformer based)	GA ^a	12.1	20.2	11.7	10.7	19.9	10.4
	RCONCAT ^a	11.9	20.1	11.5	11.0	20.4	10.5
	ARC+L2STOP* ^c	19.5	17.1	19.0	16.7	18.8	16.3
Inputs (L, V) (transformer based)	VisualBERT(8l)	10.4	21.3	10.0	9.9	21.7	9.5
	VisualBERT(4l)	14.3	17.7	13.7	11.8	18.3	11.5
	VLN Transformer(4l) ^b	12.2	18.9	12.0	12.8	20.4	11.8
	VLN Transformer(8l) ^b	13.2	19.8	12.7	13.1	21.1	12.3
	VLN Transformer(8l) + M50 + style * ^b	15.0	20.3	14.7	16.2	20.8	15.7
Inputs (L, V) + JD / HT** (non-transformer based)	ORAR (ResNet pre-final)* ^d	26.0	15.0	-	25.3	16.2	-
	ORAR (ResNet 4th-to-last)* ^d	29.9	11.1	-	29.1	11.7	-
Inputs (L, V) + Path Traces (transformer based)	VLN Transformer(8l)	11.2	23.4	10.7	11.5	23.9	10.8
	VisualBERT(4l)	16.2	18.7	15.7	15.0	20.1	14.5
	FL _{PM} (4l) + VLN Transformer(8l)	29.9	23.4	26.8	28.2	23.8	25.6
	FL _{PM} (4l) + VisualBERT(4l)	33.0	23.6	29.5	33.4	23.8	29.7

Frameworks from ^a Chen et al. (2019), ^b Zhu et al. (2021), ^c Xiang et al. (2020), and ^d Schumann and Riezler (2022).

* Results reported by the authors.

** Systems receive two types of features - Junction Type and Heading Delta - as inputs.

Table 1: Performance on the Touchdown benchmark ranked by TC on the test partition. Systems are grouped by input types during VLN and the use of transformer blocks in architectures. Contributions of the FL_{PM} framework and path traces to improved performance are demonstrated with results for systems with two baseline transformer-based architectures - VisualBERT and VLN Transformer. These baselines also are assessed in two sizes to test the benefits of adding transformer blocks.

4 EXPERIMENTS

Our starting point in evaluating the PM-VN module and FL_{PM} is performance in relation to benchmark systems (see Table 1). Ablations are conducted by removing individual operations (see Table 2) and the role of training data is assessed (see Table 3). To minimise computational cost, we implement frameworks with low numbers of layers and attention heads in transformer models.

4.1 Experiment Settings

Metrics We align with Chen et al. (2019) in reporting task completion (TC), shortest-path distance (SPD), and success weighted edit distance (SED) for ϕ_{VLN} . All metrics are derived using the Touchdown navigation graph. TC is a binary measure of success 0, 1 in ending a route with a prediction $c_{t-1}^o = y_{t-1}^o$ or $c_{t-1}^o = y_{t-1}^{o-1}$ and SPD is calculated as the mean distance between c_{t-1}^o and y_{t-1}^o . SED is the Levenshtein distance between the predicted path in relation to the defined route path and is only applied when TC = 1.

Hyperparameter Settings Frameworks are trained for 80 epochs with batch size=30. Scores are reported for the epoch with the highest SPD on $\mathcal{D}_{\phi_{VLN}}^{Dev}$. Pretraining for the PM-VLN module is conducted for 10 epochs with batch sizes $\phi_1 = 60$ and $\phi_2 = 30$. Frameworks are optimised using AdamW with a learning rate of 2.5×10^{-3} (Loshchilov and Hutter, 2017).

4.2 Touchdown

Experiment Design: Chen et al. (2019) define two separate tasks in the Touchdown benchmark: VLN and spatial description resolution. This research aligns with other studies (Zhu et al., 2021, 2022) in conducting evaluation on the navigation component as a standalone task. **Dataset and Data Preprocessing:** Frameworks are evaluated on full partitions of Touchdown with $D^{Train} = 6,525$, $D^{Dev} = 1,391$, and $D^{Test} = 1,409$ routes. Trajectory lengths vary with $D^{Train} = 34.2$, $D^{Dev} = 34.1$, and $D^{Test} = 34.4$ mean steps per route. Junction Type and Heading Delta are additional inputs generated from the environment graph and low-level visual features (Schumann and Riezler, 2022). M-50 + style is a subset of the StreetLearn dataset with $D^{Train} = 30,968$ routes of 50 nodes or less and multimodal style transfer applied to instructions (Zhu et al., 2021). **Embeddings:** All architectures evaluated in this research receive the same base cross-modal embeddings x_{η} proposed by Zhu et al. (2021), which are learned by a combination of the outputs of a pretrained BERT-base encoder with 12 encoder layers and attention heads. At each step, a fully connected layer is used for textual embeddings ζ_t and a 3 layer CNN returns the perspective ψ_t . FL_{PM} frameworks also receive an embedding of the path trace tr_t at step t . As this constitutes additional signal on the route, we evaluate a VisualBERT model (4l) that also receives tr_t , which in this case is combined with ψ_t ahead of inclusion in x_{η_t} . **Results:** In Table 1 the first block of frameworks consists of architectures composed primarily of

convolutional and recurrent layers. VLN Transformer is a framework proposed by Zhu et al. (2021) for the Touchdown benchmark and consists of a transformer-based cross-modal encoder with 8 encoder layers and 8 attention heads. VLN Transformer + M50 + style is a version of this framework pretrained on the dataset described above. To our knowledge, this was the transformer-based framework with the highest TC on Touchdown preceding our work. ORAR (ResNet 4th-to-last) (Schumann and Riezler, 2022) is from work published shortly before the completion of this research and uses two types of features to attain highest TC in prior work. Standalone VisualBERT models are evaluated in two versions with 4 and 8 layers and attention heads. A stronger performance by the smaller version indicates that adding self-attention layers is unlikely to improve VLN predictions. This is further supported by the closely matched results for the VLN Transformer(4l) and VLN Transformer(8l). FL_{PM} frameworks incorporate the PM-VLN module pretrained on auxiliary tasks (ϕ_1, ϕ_2) - and one of VisualBERT (4l) or VLN Transformer(8l) as the main model. Performance on TC for both of these transformer models doubles when integrated into the framework. A comparison of results for standalone VisualBERT and VLN Transformer systems with path traces supports the use of specific architectural components that can exploit this additional input type. Lower SPD for systems run with the FL_{PM} framework reflect a higher number of routes where a stop action was predicted prior to route completion. Although not a focus for the current research, this shortcoming in VLN benchmarks has been addressed in other work (Xiang et al., 2020; Blukis et al., 2018).

4.3 Assessment of Specific Operations

Ablations are conducted on the framework with the highest TC i.e. FL_{PM} + VisualBERT(4l). The tests do not provide a direct measure of operations as subsequent computations in forward and backward passes by retained components are not accounted for. Results indicate that initial alignment is critical to cross-modal prioritisation and support the use of in-domain data during pretraining.

	Development		
	TC \uparrow	SPD \downarrow	SED \uparrow
FL _{PM} + VisualBERT(4l)	33.0	23.6	29.5
PM-VLN			
- $g_{PMT P}$ (a)	7.1	26.8	6.8
- g_{PMF} minus g_{VBF} (b)	27.9	25.7	24.9
- g_{PMF} minus t_{t-1} (c)	29.8	21.8	27.2
FL _{PM}			
- g_{Attn} with g_{Cat} (d)	18.8	30.5	16.4
- $g_{Clas_{max\ x_i}}$ with g_{Clas} (e)	31.7	21.9	28.2

Table 2: Ablations on core operations in the PM-VLN (variants (a)-(c)) and the FL_{PM} framework (variants (d) and (e)).

Ablation 1: PM-VLN Prioritisation in the PM-VLN module constitutes a sequential chain of operations. Table 2 reports results for variants of the framework where the PM-VLN excludes individual operations. Starting with $g_{PMT P}$, trajectory estimation is replaced with a fixed count of 34 steps for each route tr_t (see variant (a)). This deprives the PM-VLN of a method to take account of the current route when synchronising τ and sequences of visual inputs. All subsequent operations are impacted and the variant reports low scores for all metrics. Two experiments are then conducted on g_{PMF} . In variant (b), visual boost filtering is disabled and feature-level localisation relies on a base ψ_t . A variant excluding linguistic components from g_{PMF} is then implemented by specifying t_t as the default input from τ_t (see variant (c)). In practice, span selection in this case is based on trajectory estimation only.

Ablation 2: FL_{PM} Ablations conclude with variants of FL_{PM} where core functions are excluded from other submodules in the framework. Results for variant (d) demonstrate the impact of replacing the operation defined in Equation 3 with a simple concatenation on outputs from PM-VLN e_l and e'_v . A final experiment compares methods for generating action predictions: in variant (e), $g_{Clas_{max\ x_i}}$ is replaced by the standard implementation for classification in VisualBERT. Classification with dropout and a single linear layer underperforms our proposal by 1.3 points on TC.

4.4 Assessment of Training Strategy

A final set of experiments is conducted to measure the impact of training data for auxiliary tasks (ϕ_1, ϕ_2).

Training Strategy 1: Exploiting Street Pattern in Trajectory Estimation We conduct tests on alternate samples to examine the impact of route types in $D_{\phi_1}^{Train}$. The module for FL_{PM} frameworks in Table 1 is trained on path traces drawn from an area in central Pittsburgh (see **SupMat:Sec.3**) with a rectangular street pattern that aligns with the urban grid type Lynch (1981) found in the location of routes in Touchdown. Table 3 presents results for modules trained on routes selected at random outside of this area. In variants (f) and (g), versions V2 and V3 of $D_{\phi_1}^{Train}$ each consist of 17,000 samples drawn at random from the remainder of a total set of 70,000 routes. Routes that conform to curvilinear grid types are observable in outer areas of Pittsburgh. Lower TC for these variants prompts consideration of street patterns when generating path traces. A variant (h) where the $g_{PMT P}$ submodule receives no pretraining underlines - along with variant (a) in Table 2 - the importance of the initial alignment step to our proposed method of cross-modal prioritisation.

Training Strategy 2: In-domain Data and Feature-Level Localisation We conclude by examining the use of in-domain data when pretraining the g_{PMF} submodule ahead of feature-level localisation operations in the PM-VLN. In Table 3, versions of FL_{PM} are evaluated subsequent to pretraining with varying sized subsets of the Conceptual Captions dataset Sharma et al. (2018). This resource of general image-text pairs is selected as it has been proposed for pretraining VLN systems (see below). Samples

are selected at random and grouped into two training partitions equivalent in number to 100% (variant(i)) and 150% of $D_{\phi_2}^{Train}$ (variant (j)). In place of the multi-objective loss applied to the MC-10 dataset, $\theta_{g_{PMF}}$ are optimised on a single goal of cross-modal matching. Variant (k) assesses FL_{PM} when no pretraining for the g_{PMF} submodule is undertaken. Lower results for variants (i), (j), and (k) support pretraining on small sets of in-domain data as an alternative to optimising VLN systems on large-scale datasets of general samples.

	Development		
	TC \uparrow	SPD \downarrow	SED \uparrow
FL_{PM} + VisualBERT(4l)	33.0	23.6	29.5
Pretraining for $g_{PMT P}$			
- $g_{PMT P} + D_{\phi_1}^{Train}V2$ (f)	11.9	20.1	11.5
- $g_{PMT P} + D_{\phi_1}^{Train}V3$ (g)	13.6	20.5	13.1
- $g_{PMT P}$ no pretraining (h)	4.7	27.6	1.9
Pretraining for g_{PMF}			
- $g_{PMF} + D_{\phi_2}^{Train}V2$ (i)	19.8	23.2	17.2
- $g_{PMF} + D_{\phi_2}^{Train}V3$ (j)	23.9	20.8	20.3
- g_{PMF} no pretraining (k)	6.3	25.1	4.6

Table 3: Assessment of the pretraining strategy for individual PM-VLN submodules $g_{PMT P}$ (variants (f) to (h)) and g_{PMF} (variants (i) to (k)) using alternative datasets for auxiliary tasks. Variants are also run with no pretraining of $g_{PMT P}$ and g_{PMF} .

5 RELATED WORK

This research aims to extend cross-disciplinary links between machine learning and computational cognitive neuroscience in the study of prioritisation in attention. This section starts with a summary of literature in these two disciplines that use computational methods to explore this subject. Our training strategy is positioned in the context of prior work on pretraining frameworks for VLN. The section concludes with work related to the alignment and feature-level operations performed by our PM-VLN model.

Computational Implementations of Prioritisation in Attention Denil et al. (2012) proposed a model that generates saliency maps where feature selection is dependent on high-level signals in the task. The full system was evaluated on computer vision tasks where the aim is to track targets in video. A priority map computation was implemented in object detection models by Wei et al. (2016) to compare functions in these systems to those observed in human visual attention. Anantrasirichai et al. (2017) used a Support Vector Machine classifier to model visual attention in human participants traversing four terrains. Priority maps were then generated to study the interaction of prioritised features and a high-level goal of maintaining smooth locomotion. A priority map component was incorporated into a CNN-based model of primate attention mechanisms by Zelinsky and Adeli (2019) to prioritise locations containing classes of interest when

performing visual search. Studies on spatial attention in human participants have explored priority map mechanisms that process inputs consisting of auditory stimuli and combined linguistic and visual information (Golob et al., 2017; Cavicchio et al., 2014). To our knowledge, our work is the first to extend neuropsychological work on prioritisation over multiple modalities to a computational implementation of a cross-modal priority map for machine learning tasks.

Pretraining for VLN Tasks Two forms of data samples - in-domain and generic - are used in pretraining prior to conducting VLN tasks. In-domain data samples have been sourced from image-caption pairs from online rental listings (Guhur et al., 2021) and other VLN tasks (Zhu et al., 2021). In-domain samples have also been generated by augmenting or reusing in-task data (Fried et al., 2018; Huang et al., 2019; Hao et al., 2020; He et al., 2021; Schumann and Riezler, 2022). Generic samples from large-scale datasets designed for other Vision-Language tasks have been sourced to improve generalisation in transformer-based VLN agents. Majumdar et al. (2020) conduct large-scale pretraining with 3.3M image-text pairs from Conceptual Captions Sharma et al. (2018) and Qi et al. (2021) initialise a framework with weights trained on four out-of-domain datasets. Our training strategy relies on datasets with a few thousand samples derived from resources where additional samples are available at low cost.

Methods for Aligning and Localising Features in Linguistic and Visual Sequences Alignment in multimodal tasks is often posited as an implicit subprocess in an attention-based component of a transformer (Tsai et al., 2019; Zhu et al., 2021). Huang et al. (2019) identified explicit cross-modal alignment as an auxiliary task that improves agent performance in VLN. Alignment in this case is measured as a similarity score on inputs from the main task. In contrast, our PM-VLN module conducts a hierarchical process of trajectory planning and learned localisation to pair inputs. A similarity measure was the basis for an alignment step in the Vision-Language Pretraining framework proposed by Li et al. (2021). A fundamental point of difference with our work is that this framework - along with related methods (Jia et al., 2021) - is trained on a distinct class of tasks where the visual input is a single image as opposed to a temporal sequence. Several VLN frameworks containing components that perform feature localisation on visual inputs have been pretrained on object detection (Majumdar et al., 2020; Suglia et al., 2021; Hu et al., 2019). In contrast, we include visual boost filtering in g_{PMF} to prioritise visual features. Our method of localising spans using a concatenation of the enhanced visual input and cross-modal embeddings is unique to this research.

6 CONCLUSION

We take inspiration from a mechanism described in neurophysiological research with the introduction of a priority map module that combines temporal sequence alignment enabled by high-level trajectory estimation and feature-level localisation. Two new resources comprised of in-domain samples and a tailored training strategy are proposed to enable data-efficient pretraining of the

PM-VLN module ahead of the main VLN task. A novel framework enables action prediction with maxout activation on a combination of the outputs from the PM-VLN module and a transformer-based encoder. Evaluations demonstrate that our module, framework, and pretraining strategy double the performance of standalone transformers in outdoor VLN.

7 ACKNOWLEDGMENTS

This publication is supported by the Digital Visual Studies program at the University of Zurich and funded by the Max Planck Society. RS has received funding by the Swiss National Science Foundation (project MUTAMUR; no. 176727). The authors would like to thank Howard Chen, Piotr Mirowski, and Wanrong Zhu for assistance and feedback on questions related to the Touchdown task, data, and framework evaluations.

REFERENCES

- N. Anantrasirchai, K. A. Daniels, J. F. Burn, I. D. Gilchrist, and D. R. Bull. Fixation prediction and visual priority maps for biped locomotion. *IEEE Transactions on Cybernetics*, 48(8):2294–2306, 2017.
- J. Armitage, E. Kacupaj, G. Tahmasebzadeh, M. Maleshkova, R. Ewerth, and J. Lehmann. Mlm: A benchmark dataset for multitask learning with multiple languages and modalities. In *Proceedings of the 29th ACM International Conference on Information & Knowledge Management*, pages 2967–2974, 2020.
- B. Bayer. A method for the digital enhancement of unsharp, grainy photographic images. *Advances in Computer Vision and Image Processing*, 2:Chapter–2, 1986.
- J. W. Bisley and K. Mirpour. The neural instantiation of a priority map. *Current Opinion in Psychology*, 29:108–112, 2019.
- V. Blukis, D. Misra, R. A. Knepper, and Y. Artzi. Mapping navigation instructions to continuous control actions with position-visitation prediction. In *Conference on Robot Learning*, pages 505–518. PMLR, 2018.
- J. Carranza-Rojas, S. Calderon-Ramirez, A. Mora-Fallas, M. Granados-Menani, and J. Torrents-Barrena. Unsharp masking layer: injecting prior knowledge in convolutional networks for image classification. In *International Conference on Artificial Neural Networks*, pages 3–16. Springer, 2019.
- F. Cavicchio, D. Melcher, and M. Poesio. The effect of linguistic and visual salience in visual world studies. *Frontiers in Psychology*, 5:176, 2014.
- H. Chen, A. Suhr, D. Misra, N. Snively, and Y. Artzi. Touchdown: Natural language navigation and spatial reasoning in visual street environments. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12538–12547, 2019.
- S. Chen, Y. Zhao, Q. Jin, and Q. Wu. Fine-grained video-text retrieval with hierarchical graph reasoning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10638–10647, 2020.
- M. Denil, L. Bazzani, H. Larochelle, and N. de Freitas. Learning where to attend with deep architectures for image tracking. *Neural Computation*, 24(8):2151–2184, 2012.
- J. H. Fecteau and D. P. Munoz. Saliency, relevance, and firing: a priority map for target selection. *Trends in Cognitive Sciences*, 10(8):382–390, 2006.
- D. Fried, R. Hu, V. Cirik, A. Rohrbach, J. Andreas, L.-P. Morency, T. Berg-Kirkpatrick, K. Saenko, D. Klein, and T. Darrell. Speaker-follower models for vision-and-language navigation. *Advances in Neural Information Processing Systems*, 31, 2018.
- E. J. Golob, K. B. Venable, J. Scheuerman, and M. T. Anderson. Computational modeling of auditory spatial attention. In *CogSci*, 2017.
- I. Goodfellow, D. Warde-Farley, M. Mirza, A. Courville, and Y. Bengio. Maxout networks. In *International Conference on Machine Learning*, pages 1319–1327. PMLR, 2013.
- J. Gottlieb, M. Cohanpour, Y. Li, N. Singletary, and E. Zabe. Curiosity, information demand and attentional priority. *Current Opinion in Behavioral Sciences*, 35:83–91, 2020.
- P.-L. Guhur, M. Tapaswi, S. Chen, I. Laptev, and C. Schmid. Airbert: In-domain pretraining for vision-and-language navigation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 1634–1643, 2021.
- W. Hao, C. Li, X. Li, L. Carin, and J. Gao. Towards learning a generic agent for vision-and-language navigation via pre-training. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13137–13146, 2020.
- B. Y. Hayden and J. L. Gallant. Combined effects of spatial and feature-based attention on responses of v4 neurons. *Vision Research*, 49(10):1182–1187, 2009.
- K. He, Y. Huang, Q. Wu, J. Yang, D. An, S. Sima, and L. Wang. Landmark-rxr: Solving vision-and-language navigation with fine-grained alignment supervision. *Advances in Neural Information Processing Systems*, 34, 2021.
- K. M. Hermann, M. Malinowski, P. Mirowski, A. Banki-Horvath, K. Anderson, and R. Hadsell. Learning to follow directions in street view. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 11773–11781, 2020.
- R. Hu, D. Fried, A. Rohrbach, D. Klein, T. Darrell, and K. Saenko. Are you looking? grounding to multiple modalities in vision-and-language navigation. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 6551–6557. Florence, Italy, July 2019. Association for Computational Linguistics. doi: 10.18653/v1/P19-1655. URL <https://aclanthology.org/P19-1655>.
- H. Huang, V. Jain, H. Mehta, A. Ku, G. Magalhaes, J. Baldrige, and E. Ie. Transferable representation learning in vision-and-language navigation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 7404–7413, 2019.
- L. Itti and C. Koch. A saliency-based search mechanism for overt and covert shifts of visual attention. *Vision Research*, 40(10-12):1489–1506, 2000.
- C. Jia, Y. Yang, Y. Xia, Y.-T. Chen, Z. Parekh, H. Pham, Q. Le, Y.-H. Sung, Z. Li, and T. Duerig. Scaling up visual and vision-language representation learning with noisy text supervision. In *International Conference on Machine Learning*, pages 4904–4916. PMLR, 2021.
- T. Le, K. Pho, T. Bui, H. T. Nguyen, and M. Le Nguyen. Object-less vision-language model on visual question classification for blind people. In *ICAAIT (3)*, pages 180–187, 2022.
- J. Li, R. Selvaraju, A. Gotmare, S. Joty, C. Xiong, and S. C. H. Hoi. Align before fuse: Vision and language representation learning with momentum distillation. *Advances in Neural Information Processing Systems*, 34, 2021.
- L. H. Li, M. Yatskar, D. Yin, C.-J. Hsieh, and K.-W. Chang. Visualbert: A simple and performant baseline for vision and language. *arXiv preprint arXiv:1908.03557*, 2019.
- X. Li, X. Yin, C. Li, P. Zhang, X. Hu, L. Zhang, L. Wang, H. Hu, L. Dong, F. Wei, et al. Oscar: Object-semantics aligned pre-training for vision-language tasks. In *European Conference on Computer Vision*, pages 121–137. Springer, 2020.
- Y.-B. Lin and Y.-C. F. Wang. Audiovisual transformer with instance attention for audio-visual event localization. In *Proceedings of the Asian Conference on Computer Vision*, 2020.
- Z. Liu, H. Mao, C.-Y. Wu, C. Feichtenhofer, T. Darrell, and S. Xie. A convnet for the 2020s. *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022.
- I. Loshchilov and F. Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017.
- K. Lynch. A theory of good city form. ma. *Cambridge: Massachusetts Institute of Technology*, 1981.
- A. Majumdar, A. Shrivastava, S. Lee, P. Anderson, D. Parikh, and D. Batra. Improving vision-and-language navigation with image-text pairs from the web. In *European Conference on Computer Vision*, pages 259–274. Springer, 2020.
- P. Mirowski, M. Grimes, M. Malinowski, K. M. Hermann, K. Anderson, D. Teplyashin, K. Simonyan, A. Zisserman, R. Hadsell, et al. Learning to navigate in cities without a map. *Advances in Neural Information Processing Systems*, 31:2419–2430, 2018.

- R. Ptak. The frontoparietal attention network of the human brain: action, saliency, and a priority map of the environment. *The Neuroscientist*, 18(5):502–515, 2012.
- Y. Qi, Z. Pan, Y. Hong, M.-H. Yang, A. van den Hengel, and Q. Wu. The road to know-where: An object-and-room informed sequential bert for indoor vision-language navigation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 1655–1664, 2021.
- R. Schumann and S. Riezler. Analyzing generalization of vision and language navigation to unseen outdoor areas. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 7519–7532, Dublin, Ireland, May 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.acl-long.518. URL <https://aclanthology.org/2022.acl-long.518>.
- P. Sharma, N. Ding, S. Goodman, and R. Soricut. Conceptual captions: A cleaned, hypernymed, image alt-text dataset for automatic image captioning. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2556–2565, 2018.
- H. Shinoda, M. M. Hayhoe, and A. Shrivastava. What controls attention in natural environments? *Vision Research*, 41(25-26):3535–3545, 2001.
- W. Su, X. Zhu, Y. Cao, B. Li, L. Lu, F. Wei, and J. Dai. VL-BERT: pre-training of generic visual-linguistic representations. In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net, 2020.
- A. Suglia, Q. Gao, J. Thomason, G. Thattai, and G. Sukhatme. Embodied bert: A transformer model for embodied, language-guided visual task completion. *arXiv preprint arXiv:2108.04927*, 2021.
- B. W. Tatler, M. M. Hayhoe, M. F. Land, and D. H. Ballard. Eye guidance in natural vision: Reinterpreting saliency. *Journal of Vision*, 11(5):5–5, 2011.
- Y.-H. H. Tsai, S. Bai, P. P. Liang, J. Z. Kolter, L.-P. Morency, and R. Salakhutdinov. Multimodal transformer for unaligned multimodal language sequences. In *Proceedings of the conference. Association for Computational Linguistics. Meeting*, volume 2019, page 6558. NIH Public Access, 2019.
- Z. Wei, H. Adeli, M. H. Nguyen, G. Zelinsky, and D. Samarasinghe. Learned region sparsity and diversity also predicts visual attention. *Advances in Neural Information Processing Systems*, 29, 2016.
- J. Xiang, X. Wang, and W. Y. Wang. Learning to stop: A simple yet effective approach to urban vision-language navigation. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 699–707, Online, Nov. 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.findings-emnlp.62. URL <https://aclanthology.org/2020.findings-emnlp.62>.
- G. J. Zelinsky and H. Adeli. Learning to attend in a brain-inspired deep neural network. *Journal of Vision*, 19(10):282d–282d, 2019.
- G. J. Zelinsky and J. W. Bisley. The what, where, and why of priority maps and their interactions with visual working memory. *Annals of the New York Academy of Sciences*, 1339(1):154, 2015.
- W. Zhu, X. Wang, T.-J. Fu, A. Yan, P. Narayana, K. Sone, S. Basu, and W. Y. Wang. Multimodal text style transfer for outdoor vision-and-language navigation. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 1207–1221, 2021.
- W. Zhu, Y. Qi, P. Narayana, K. Sone, S. Basu, X. Wang, Q. Wu, M. Eckstein, and W. Y. Wang. Diagnosing vision-and-language navigation: What really matters. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 5981–5993, Seattle, United States, July 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.naacl-main.438. URL <https://aclanthology.org/2022.naacl-main.438>.

Appendices

A NOTATION

Notations used in multiple sections of this paper are defined here for fast reference. Auxiliary tasks (ϕ_1, ϕ_2) and the main VLN task ϕ_{VLN} constitute the set of tasks Φ . Inputs and embeddings are specified as l (linguistic), v (visual), and η (multimodal). A complete textual instruction is denoted as τ , ζ is a span, and ψ is a perspective. Linguistic and visual inputs for the PM-VLN are denoted as (i'_t, ψ_t) and embeddings processed in prioritisation operations are $(e_l, e_v)_t$. In contrast, U denotes a set of embeddings from the main model, which are derived from inputs $(\bar{e}_\eta, \psi_{cat})$. The notations Δ and \oplus are respectively visual boost filtering and self-attention operations. Table 4 provides a reference source for standard notation appearing throughout this paper. Other notations are defined in the sections where they are used.

Notation	Usage in this paper
A	Matrix
AA	Identity matrix
B, b	Bias
\mathcal{D}	Dataset
$Train, Dev, Test$	Dataset partitions
\exists	Exists
\forall	For every (eg member in a set)
g	Function
H	Hypothesis
\mathcal{L}	Layer of a model
len	Length
μ	Mean
n	Number of samples
\vee	Or
P	Probability
q	Algorithm
S	Signal detected
σ	Standard deviation
Θ	Set of parameters
W, w	Set of weights
$ x $	Sequence
\triangleq	Equal by definition

Table 4: Reference List for Standard Notation.

B DATASETS

B.1 Generation and Partition Sizes

The MC-10 dataset consists of visual, textual and geospatial data for landmarks in 10 US cities. We generate the dataset with a modified version of the process outlined by Armitage et al. (2020). Two base entity IDs - Q2221906 (“geographic location”) and Q83620 (“thoroughfare”) - form the basis of queries to extract entities at a distance of ≤ 2 hops in the Wikidata knowledge graph¹. Constituent cities consist of incorporated places exceeding 1 million people ranked by population density based on data for April 1, 2020 from the US Census Bureau². Images and coordinates are sourced from Wikimedia and text summaries are extracted with the MediaWiki API. Geographical cells are generated using the S2 Geometry Library³ with a range of n entities [1, 5]. Statistics for MC-10 are presented by partition in Table 5. As noted above, only a portion of textual inputs are used in pretraining and experiments.

¹ <https://query.wikidata.org/>

² <https://www.census.gov/programs-surveys/decennial-census/data/datasets.html>

³ <https://code.google.com/archive/p/s2-geometry-library/>

C CODE AND DATA

Source code for the project and instructions to run the framework are released and maintained in a public GitHub repository under MIT license (<https://github.com/JasonArmitage-res/PM-VLN>). Code for the environment, navigation, and training adheres to the codebases released by Zhu et al. (2021) and Chen et al. (2019) with the aim of enabling comparisons with benchmarks introduced in prior work on Touchdown. Full versions of the MC-10 and TR-NY-PIT-central datasets are published on Zenodo under Creative Commons public license (<https://zenodo.org/record/6891965#.YtwoS3ZBxD8>).