

Information Fusion in Visual-Task Inference

Amin Haji-Abolhassani
Centre for Intelligent Machines
McGill University
Montreal, Quebec H3A 2A7, Canada
Email: amin@cim.mcgill.ca

James J. Clark
Centre for Intelligent Machines
McGill University
Montreal, Quebec H3A 2A7, Canada
Email: clark@cim.mcgill.ca

1 **Abstract**—Eye movement is a rich modality that can provide
2 us with a window into a person’s mind. In a typical human-
3 human interaction, we can get information about the behavioral
4 state of the others by examining their eye movements. For
5 instance, when a poker player looks into the eyes of his
6 opponent, he looks for any indication of bluffing by verifying
7 the dynamics of the eye movements. However, the information
8 extracted from the eyes is not the only source of information
9 we get in a human-human interaction and other modalities,
10 such as speech or gesture, help us infer the behavioral state of
11 the others. Most of the time this fusion of information refines
12 our decisions and helps us better infer people’s cognitive and
13 behavioral activity based on their actions. In this paper, we
14 develop a probabilistic framework to fuse different sources of
15 information to infer the ongoing task in a visual search activity
16 given the viewer’s eye movement data. We propose to use a
17 dynamic programming method called *token passing* in an eye-
18 typing application to reveal what the subject is typing during
19 a search process by observing his direction of gaze during
20 the execution of the task. Token passing is a computationally
21 simple technique that allows us to fuse higher order constraints
22 in the inference process and build models dynamically so we
23 can have unlimited number of hypotheses. In the experiments
24 we examine the effect of higher order information, in the form
25 of a lexicon dictionary, on the task recognition accuracy.

26 **Keywords**—attention; visual search; cognitive modeling; task
27 inference; information fusion; eye movement;

28 I. INTRODUCTION

29 The link between eye movements and visual task has en-
30 joyed burgeoning attention in psychophysical and cognitive
31 sciences. Particularly the effect of visual task on parameters
32 of eye movements have been investigated for a long time
33 in the literature. In two seminal studies, Yarbus [1967] and
34 Buswell [1935] showed that visual task has a great influence
35 on specific parameters of eye trajectory. Figure 1 shows
36 Yarbus’s observation that implies fixation locations are not
37 randomly distributed in a scene but instead tend to cluster
38 on some regions at the expense of others. In this figure we
39 see how visual task can modulate the conspicuity of different
40 regions and as a result change the pattern of eye movements.
41 Based on this experiment, the effect of task on the pattern
42 an parameters of eye movement is called the *forward Yarbus*
43 *process*.

44 The forward Yarbus process is also studied in a work
45 by Clark and O’Regan [1998], who examined the dynamics

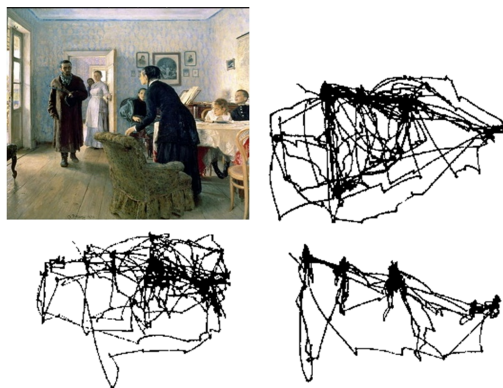


Figure 1: Eye trajectories recorded by Yarbus while a viewer carried out different visual tasks. Upper right - no specific task, lower left - estimate the wealth of the family, lower right - give the ages of the people in the painting [Yarbus, 1967].

46 of eye movements when reading a text. They showed that
47 when reading a text the centre of gaze (COG) lands on the
48 locations that minimize the ambiguity of the word arising
49 from the incomplete recognition of the letters. In another
50 study Castelhana et al. [2009] showed that different patterns
51 of eye movements emerge from tasks of memorizing a scene
52 and searching for an object in it.

53 In a forward Yarbus process the visual task is given as
54 an input and the output is task-dependent scanpaths of eye
55 movements. The other way of looking at the interaction
56 between eye movements and visual task is to study the
57 reverse path from the scanpaths to the visual task. The
58 inference of task from eye movement data is called an
59 *inverse Yarbus process* and has recently gained a growing
60 interest in psychophysical studies of human attention.

61 Visual search is one of the main ingredients of human
62 vision that plays an important role in our everyday life.
63 Recently in [Haji-Abolhassani and Clark, 2011a,b] we pro-
64 posed a model based on the theory of Hidden Markov
65 Models (HMMs) to infer what the viewer is looking for
66 in a task of searching for pop-out objects in digitally

67 created stimuli. In [Haji-Abolhassani and Clark, 2012a,b] we
 68 extended our model to infer the word that the viewers typed
 69 using their eye movements (*eye-typing*) in a soft keyboard
 70 application. In both scenarios the viewer executes visual
 71 search and the model calculates a probability distribution
 72 on different possible tasks given the eye data, and makes
 73 an inference about what objects are being sought using
 74 *maximum likelihood* (ML). In real life, however, we incor-
 75 porate a-priori sources of information in the ML estimator
 76 and make inferences based on *maximum a-posteriori* (MAP)
 77 estimation. For instance, when looking for an orange in
 78 a basket of fruits, the prior knowledge about the color of
 79 oranges helps us skip the objects with different colors and
 80 narrow down our search to the orange areas.

81 In this paper we extend our HMM-based model presented
 82 in [Haji-Abolhassani and Clark, 2012a,b]; which we will
 83 be referring to as tri-state HMM (TSHMM) in the rest¹¹⁹
 84 of the text; to incorporate a-priori information to infer¹²⁰
 85 the visual task in the eye-typing application. In order to¹²¹
 86 infer the ongoing task, we propose to use the TSHMM¹²²
 87 model within a simple conceptual model of eye movement¹²³
 88 recognition based on a technique called *token passing* that¹²⁴
 89 incorporates the TSHMMs in a transition network structure.¹²⁵
 90 In the new structure, the higher order constraints are applied¹²⁶
 91 along transitions from a TSHMM unit to another. Moreover,¹²⁷
 92 since in token passing method the models are generated¹²⁸
 93 dynamically during the test phase, we can have an unlimited¹²⁹
 94 number of hypotheses in our experiments. ¹³⁰

95 In the following sections we will first revisit the TSHMM¹³¹
 96 model used for task inference in the eye-typing application.¹³²
 97 Then we show how we can equip the model with high level¹³³
 98 constraints. In the experiments we show how using a-priori¹³⁴
 99 information in the form of a lexicon dictionary improves the¹³⁵
 100 recognition rate. ¹³⁶

101 II. TASK INFERENCE USING HIDDEN MARKOV MODELS ¹³⁸

102 The application we designed for task inference in visual¹⁴⁰
 103 search is an eye-typing application, where subjects can type¹⁴¹
 104 a word by directing their gaze on its comprising characters.¹⁴²
 105 Figure 2 shows the schematic of the on-screen keyboard¹⁴³
 106 used in the experiments. We removed the letter “Z” from¹⁴⁴
 107 the keyboard to obtain a square layout to reduce directional¹⁴⁵
 108 bias. In order to impose visual search, we randomized the¹⁴⁶
 109 location of characters to eliminate any memory effect. ¹⁴⁷

110 Although visual attention and direction of gaze are some-¹⁴⁸
 111 times assumed to be the same, in oculomotor studies of¹⁴⁹
 112 human vision it is shown that the focus of attention (FOA)¹⁵⁰
 113 can be well away from the center of gaze (COG) [Fischer¹⁵¹
 114 and Weber, 1993]. Based on the alignment of the COG to¹⁵²
 115 the FOA we have two types of visual attention; that are¹⁵³
 116 covert and overt attention. In overt visual attention the FOA¹⁵⁴
 117 is aligned to the COG and in the covert visual attention the¹⁵⁵
 118 FOA is away from COG. ¹⁵⁶

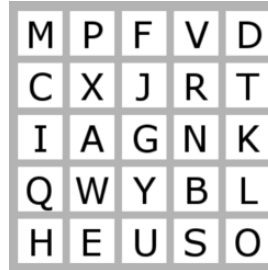


Figure 2: The schematic of the on-screen keyboard used in the eye-typing experiments. We removed the letter “Z” in order to have a square layout to reduce directional bias. Also the location of each character is randomized in each layout so that the user has to search for the characters.

Perhaps the first scientist to provide an experimental demonstration of covert attention is known to be Helmholtz [1896]. In his experiment, Helmholtz briefly illuminated inside a box by lighting a spark and looked at it through two pinholes. Before the flash he attended to a particular region of his visual field without moving his eyes in that direction. He showed that only the objects in the attended area could be recognized implying attention can be away from the eye movements.

Apart from the intrinsic difference between the FOA and the COG in covert shift of attention, the focus of overt attention can also be different from the COG reported by the eye-tracker due to the noise of the recording instrument. Moreover, overshooting or undershooting of the targets can cause a mismatch between the COG and FOA, regardless of the attention type, which urges us to allow for discrepancy between these two phenomena in the attention models.

Hidden Markov Models (HMMs) are a class of generative methods that are used to classify sequential observations. A typical HMM is composed of a number of *states* that are hidden from the observer. The transitions between the states are governed by a *transition probability matrix*, A , that gives us the chances of transitions from a state to the connecting states. At each time-step, an observation is generated according to an *observation probability density function* that is assigned to the current state. The observation pdf is characterized by a set of parameters B that defines the properties of the pdf. At the beginning of each sequence, the HMM selects a starting state according to a *initial state distribution*, Π , and carries on by chooses the next states at each time-step according to A . Figure 3 shows a sample tri-state HMM with its corresponding observation pdf, transition probabilities and initial state distribution.

HMMs have been extensively used in the field of speech recognition [Rabiner, 1990], optical character recognition [Hu et al., 1996] and anomaly detection in video surveillance [Nair and Clark, 2002] before. There are usually three different problems that are addressed in the literature

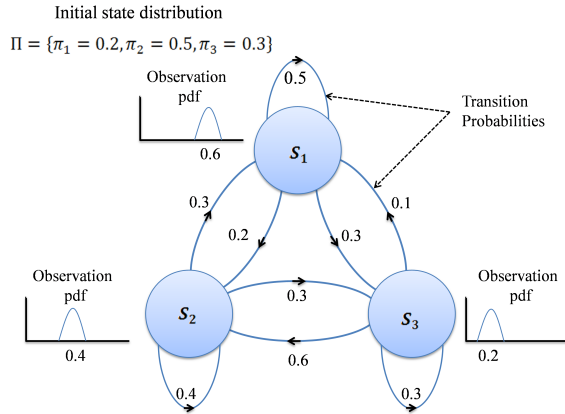


Figure 3: A sample first-order, discrete-time, continuous HMM. An HMM is defined by its number of states, transition probabilities, observation pdfs and initial state distribution in a tuple $\lambda(A, B, \Pi)$.

the lower the connecting line between clusters, the higher the similarity between the clusters). Figure 4b shows that a similar pattern appears in our experiments.

Based on these facts we designed the TSHMM model so that it allows for attentional deployment both on target and non-target objects. Furthermore, we divided the off-target fixations according to their similarity to the target. In figure 5a we see the proposed attention model for the task of looking for a character. For each character we train an HMM that has three states that represent deployment of attention on non-target, similar-to-target (*S-state*); non-target, dissimilar-to-target (*D-state*); and target characters (*T-state*). As we showed, this structure elicits more information from off-target fixations which increases the accuracy of task inference.

The observation pdfs generate COGs according to the attention state and are defined by GMMs with equal weights in the D-state and S-state, and a single Gaussian in the T-state. The GMM of the S-state has a mean vector that points to the top two similar characters according to the fixation frequency histogram (similar to figure 4b) that is obtained in the training phase. The GMM of the D-state is simply the negation of the S-state and T-state’s observation pdfs (with equal weights) which points to the whole surface of the keyboard except the target and similar-to-target locations.

The initial state distribution, determines the chances of starting from each state given an observation. Figure 5b shows how we create a word model based on the HMMs of its comprising characters. When finding a target character, we assume the transition probabilities to be proportional to the initial state distribution for the next character. Although, due to some memory effects the transition probabilities and initial state distribution might not be exactly the same, the difference seems to be negligible. Beside major reduction in training, it is only by this assumption that we are able to build a model that can accommodate unlimited number of words.

III. INFORMATION FUSION USING TOKEN PASSING

Although the experiments show that our tri-state HMM (TSHMM) can reliably be used in task inference in the eye-typing application, there are other sources of information that could be applied to the inference to improve the performance of the model. Probability distribution of task priors is a source of information that we use on a daily basis to make inferences about our observations. In our application, when the model gives us a uniform distribution over characters “V” and “U”, knowing that the proceeding character was a “Q” would help us choose “U” as the eye-typed character, because that is the character that always follows “Q” in common English words.

A similar technique is used in speech processing community to improve the results of a recognizer by applying high level constraints to the character sequences [Rabiner,

related to the HMMs; that are training, decoding and evaluation. The training is done by an algorithm called *Baum-Welch*, whereby we train the parameters of HMMs using the training data. In decoding, the best sequence of states is revealed by using a method called *Viterbi*. Finally, in the evaluation, we use a method known as *forward algorithm* to find the likelihood of an observation given the parameters of an HMM.

In the TSHMM model we used HMMs to model the cognitive process of human brain that controls the COG and FOA. In the model we represented the FOA by the hidden states of an HMM and the observations of the HMM were equivalent to the COG. The only information we observe from a human eye is the COG and the FOA is hidden from us. This is inline with the structure of HMMs, where we only see the observations and the states are hidden to the observer.

Looking for a character among other characters is a visual search task that requires a combination of features to be used to locate the target [Treisman and Gelade, 1980]. This characteristic calls for an attentive, mainly serial, limited capacity attentional deployment over a limited portion of the visual field which usually entails several fixations on distractors (non-targets) before locating a target. Moreover, during the experiments we observed a pattern in these off-target fixations that implies the FOA doesn’t randomly scan the characters to seek a target, but instead tend to verify the similar characters more often than dis-similar ones. This effect is studied before in perceptual measurement of image similarity in [Keren and Baggen, 1981, Gilmore et al., 1979]. Figure 4a shows the result of an experiment in [Gilmore et al., 1979] that categorizes the characters according to their similarity. In this figure a hierarchical clustering is used to classify characters according to their similarity (i.e.,

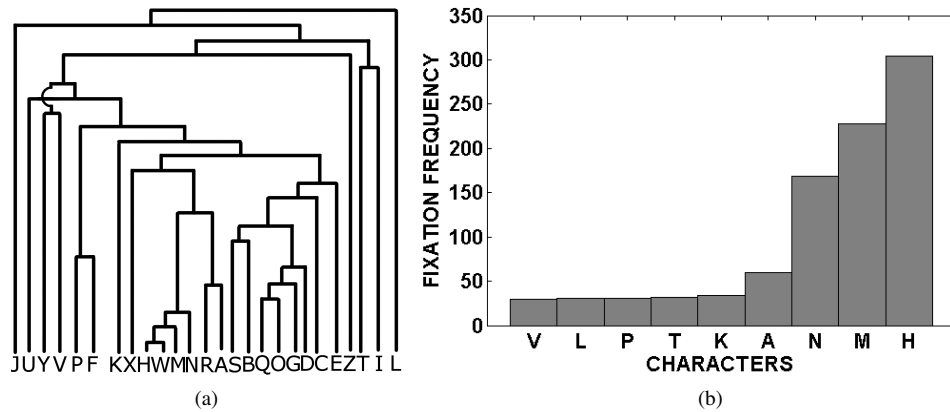


Figure 4: a) Shows the result of an experiment in perceptual measurement of image similarity that appears in [Gilmore et al., 1979, figure 1]. The results shown here are inline with what we observed in the experiments. b) Shows the top nine bars of the fixation distribution when looking for character “W”. Similar characters tend to draw attention towards themselves.

1990]. When recognizing a speech signal, the constraint is imposed to the decision making engine in the form of a lexicon dictionary called *language model* (LM) that provides us with a-prior information about the current speech part that is being pronounced given the proceeding one.

Since in our application we deal with common English words as well, we use a similar technique to apply higher order constraints on the recognizer. Depending on the order of dependency of characters, we have different orders of LM. In a *unigram* LM we assume the characters to be independent of each other. Applying a unigram LM to our TSHMM reduces the model to what we had before. However, in this paper we use a *bigram* LM whereby we impose a first order Markov process on the sequence of characters.

In order to build a LM we need to get a database of valid English words. Then we can train the bigram LM by assuming a first order Markov chain as the underlying process of character sequences. The training is done by counting the number of each pair of transitions in the corpus. In the end, a technique called *add one smoothing* is applied to the count numbers by assuming each pair occurs more than it actually does to assign non-zero probabilities to the unseen pairs in the training corpus [Huang et al., 2001, chapter 11]. Eventually the language model gives us the probability of p_{ij} for each pair of characters (i, j) , where p_{ij} is the probability of seeking character j after having found character i in our eye-typing application.

In the previous section we showed how we can train TSHMMs for each character. Therefore, by training the LM we have a complete model for the words in the dictionary that describes the transitions within the states of characters, as well as transitions between a word’s characters, in a probabilistic manner. This model can be used as a generative

model of the cognitive process of the human brain that generates eye movements during visual search for characters of a word (i.e., eye-typing). First we start from the initial state of a character according to the initial state distribution of the HMM, and by following the transition probabilities we can choose the states for each time step and generate observations according to the observation probabilities. When getting to the final state of a character, it is the language model that suggests which character, by what probability, can follow the current one.

The complete structure of the model for a two-character scenario is shown in figure 6. The bigram LM information is applied to the transitions between characters. Unlike the nodes inside the boxes that represent the states in the character models, the LM nodes don’t represent states, which means neither any observation is generated in them nor transition through them takes up any time-step in the sequence. The LM nodes are equivalent to the so-called *grammar nodes* in the speech processing literature and is merely an indication of applying LM to the model [Huang et al., 2001, page 618].

Having the generative model of eye movements during visual search, we can use the trained parameters to decode a test eye movement trajectory to infer what character sequence has been eye-typed. If we had a limited number of hypotheses (words in the dictionary), we could use *Viterbi algorithm* to classify the test data into one of the words in the dictionary [Rabiner, 1990]. Viterbi algorithm, though, requires the word models to be built beforehand to be able to compare the likelihood of each word in the dictionary. However, for a recognition task, there might be an enormous number of words in the dictionary which makes it computationally expensive to build the word model for each word statically in the Viterbi algorithm.

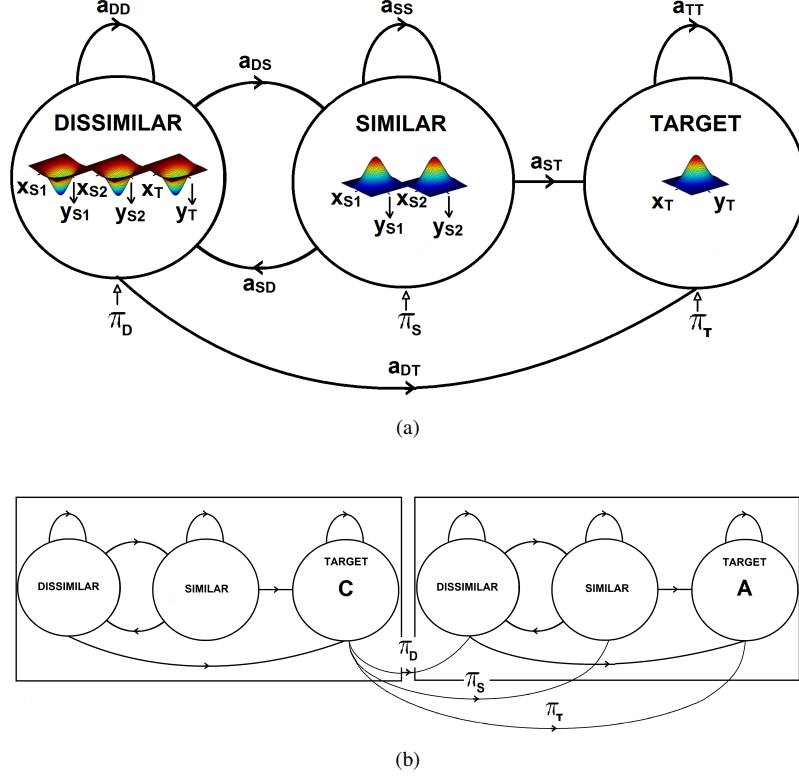


Figure 5: The structure of the tri-state HMM (TSHMM) for character recognition. a) The TSHMM for a single character. b) Concatenating the character models to build up the word “CA”. The transitions between the states are governed by the initial state probabilities.

312 An analogous problem exists in the literature related to 336
 313 speech recognition, where the dictionary of possible words 337
 314 exceeds a certain number. The technique used there, that we 338
 315 propose to be used for our problem as well, is a dynamic 339
 316 programming algorithm called *token passing* [Young et al. 340
 317 1989]. In order to find the best sequence of states that
 318 matches the observation sequence, we assign a cost R_{ij}
 319 to each transition from state i to j in an HMM equal to
 320 $\log(1/a_{ij})$ and call it the *transition cost*. a_{ij} is the transition
 321 probability in the TSHMMs model of a character if the transition is within a character model. If the transition is between characters, we use the LM statistics to evaluate the transition cost and we will have $R_{ij} = \log(1/p_{ij})$, where 341
 322 p_{ij} is the probability of going from character i to character 342
 323 j according to the LM. 343

327 The second type of cost that we use in the token passing 344
 328 method is called *local cost function* and defines the cost 345
 329 of being at state j at time t . Suppose we have data of 346
 330 the form $\langle \mathbf{Q}, y \rangle$, where $y \in Y$ is a task label in the 347
 331 set of all task labels Y and \mathbf{Q} is the vector containing the
 332 observation sequence of fixation locations $(\vec{q}_1, \vec{q}_2, \dots, \vec{q}_T)$
 333 sampled from a stochastic process $\{\vec{q}_t\}$ at discrete times
 334 $t = \{1, 2, \dots, T\}$ over random image locations denoted in
 335 Cartesian coordinates by $\vec{q}_t = (x_t, y_t)$.

The local cost function is defined as the cost of being at
 state state j at time t and is defined by $L_j(t) = \log(1/b_j(\vec{q}_t))$, where $b_j(\vec{q}_t)$ is the probability of observing \vec{q}_t at state j . Similar to the transition cost, L can be calculated using the parameters of the task-specific TSHMMs.

Having defined these two cost functions, we can calculate the alignment cost for an observation sequence Q and a sample state sequence $I = (i_0, i_1, \dots, i_T)$ by computing the alignment cost:

$$S(I) = \sum_{\tau=1:T} (R_{i_{\tau-1}i_\tau} + L_{i_\tau}(\tau)). \quad (1)$$

However, most of the time the state sequence is hidden from the observer and therefore we can't compute the alignment cost function directly. In token passing method, thus, we define a new alignment cost function called *local alignment cost function*, $s_j(t)$, that is equal to the sum of transition and local cost functions that leads to being at state j at time t .¹ Algorithm 1 shows how we can use the local

¹In Viterbi algorithm (for limited number of tasks) we can use dynamic programming to calculate the alignment cost functions using the following equation:

$$s_j(t) = \min_i [s_i(t-1) + R_{ij}] + L_j(q_t). \quad (2)$$

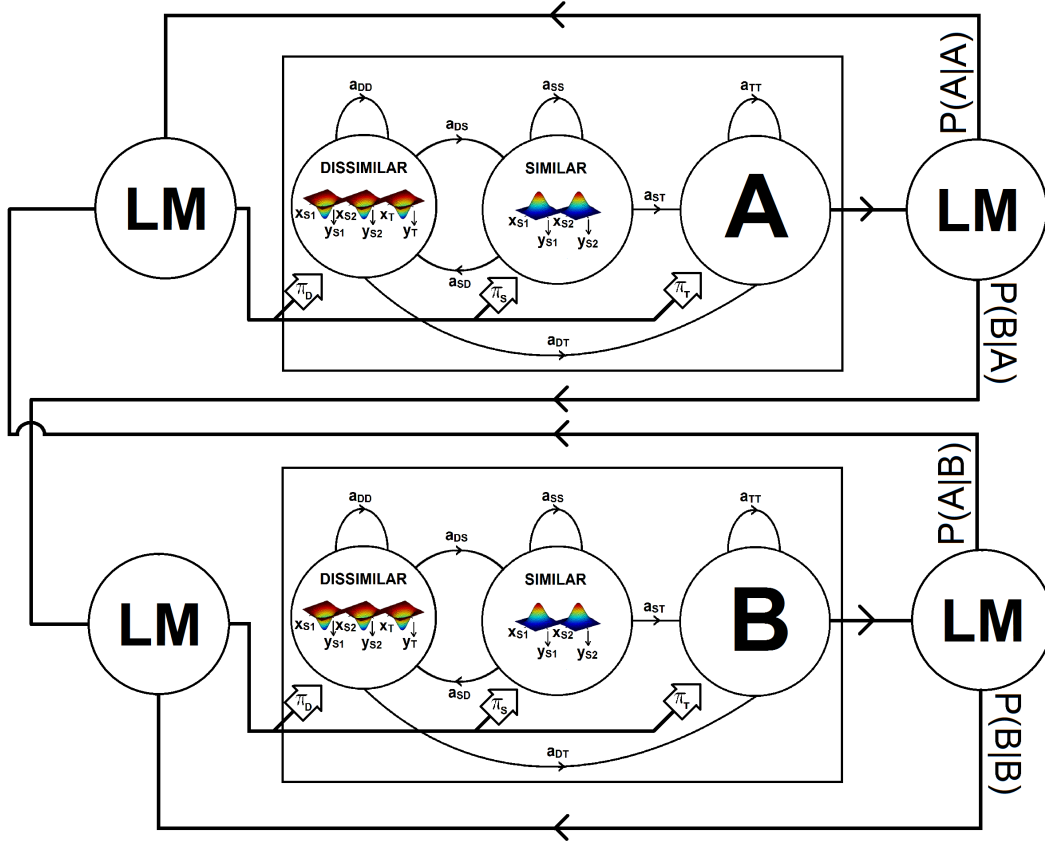


Figure 6: Incorporating the a-priori information in the form of a language model in the word model. The best state sequence for a given observation sequence can be obtained by using the token passing technique on this general word model. The LM nodes (a.k.a grammar nodes) don't generate observations, but the LM parameters are applied to the model when passing through these nodes.

348 alignment costs to decode a sequence of eye movements by₃₆₈
 349 finding the path with the minimum cost in figure 6. To do₃₆₉
 350 so, we assume that each HMM state can hold a movable₃₇₀
 351 token. We can think of a token as an object that can move₃₇₁
 352 from one state to another in our network. Each token carries₃₇₂
 353 with it a local alignment cost, which gets propagated in the₃₇₃
 354 network according to the transition and local cost functions.₃₇₄
 355 In the algorithm we refer to this cost function as the value₃₇₅
 356 of the token. At the end of the iterations, the rout with the₃₇₆
 357 minimum cost gives us the best alignment between the states₃₇₇
 358 and the observation sequence. ₃₇₈

379 IV. EXPERIMENTS ₃₈₀

360 To build a database of task-dependent eye trajectories,
 361 we ran a set of trials and recorded the eye movements of₃₈₁
 362 six subjects while eye-typing 26 different 3-character words.₃₈₂
 363 The trials started with a fixation mark of size 0.26×0.26 ₃₈₃
 364 deg appearing at the center of the screen. After foveating₃₈₄
 365 the fixation mark, the participant initiated the trial with a₃₈₅
 366 key-press. Once a trial was triggered, the word to be eye-₃₈₆
 367 typed was shown at the center of the display. Once the₃₈₇

subject indicated his readiness by pressing a key, another
 fixation mark appeared at the center followed by an on-
 screen keyboard similar to the one shown in figure 2. At
 this phase subjects eye-typed the word by searching for the
 characters appearing in it as quickly as possible and signaled
 when they were done by pressing a key (subjects were only
 told to eye-type the words as quickly as possible and press a
 key when done). Then by asking about the location of one of
 the characters (selected randomly) we verified to see if the
 subject had correctly eye-typed the words. Once the question
 is answered (by fixating the right location that contained the
 character during the experiment and pressing a bottom) the
 next word is shown and the trial carries on. ₃₈₀

The stimuli were generated by a computer and displayed
 on a 1280×800 pixel screen at a distance of 18 inches (1
 degree of visual angle corresponds to 30 pixels, approx-
 imately). Each keyboard was composed of 25 uppercase
 English characters randomly located on a 5×5 grid su-
 perimposed on a gray background (we removed the letter
 "Z" in order to have a square layout to reduce directional

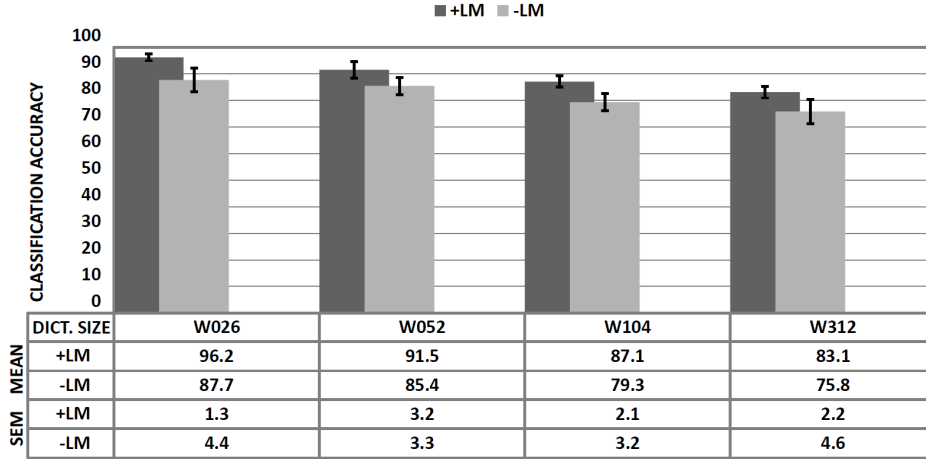


Figure 7: Comparison of the task classification accuracy using TSHMM with a bigram LM (+LM) and TSHMM with a unigram LM (-LM) in the eye-typing application. The TSHMM with a unigram LM (-LM) corresponds to previous work, where no LM was assumed. The “DICT. SIZE” row shows the number of words (hypotheses) used in each experiment with a “Wxxx” code, where “xxx” shows the number of words. Each bar shows the mean classification rate (%) of correctly recognizing the intended word in the eye-typing application. The mean value and the standard error of the mean (SEM) are represented by bars and the numerical values are given in the following table.

Algorithm 1 Token Passing algorithm

Initialize:

Assign a zero valued token to the initial states of the models.

Assign an ∞ valued token to all other states.

Algorithm:

for $t = 1$ to T **do**

for each state i **do**

 Copy the token in each state i to the connecting state j and increment its value by $R_{ij} + L_j(t)$

end for

 Discard the original tokens.

for each state i **do**

 Keep the token with the minimum value and discard the rest.

end for

end for

Termination:

The token with the minimum s value in all possible final states corresponds to the best match.

the participant’s left eye positions at 60 Hz and a chin rest was used to minimize head movements. The eye tracker’s vertical resolution is approximately 0.11 degrees and its horizontal resolution is 0.06 degrees. An LCD monitor was used for displaying the images and the subjects used both eyes to conduct the experiments.

After recording eye movements, data analysis was carried out on each trial wherein we removed the blinks, outliers and trials with wrong answers in the verification phase from the data and classified the eye movement data into saccades and fixations. Moreover, in some of the initial trials, after eye-typing the word, the viewer returned to the locations of the characters to double-check the coordinates of them. In order to simulate a real eye-typing application we removed these parts from the trajectories in the pre-processing as well.

After the preprocessing we obtained a database of 145 trajectories of the form $(\vec{q}_1, \dots, \vec{q}_T)$, each containing observation sequences of coordinates of fixations while performing the eye-typing, where $\vec{q}_t = (x_t, y_t)$ represents x -coordinate and y -coordinate of the t^{th} fixation, respectively.

In order to perform the evaluation, we compare the results of our proposed model that uses a TSHMM and a bigram LM to model the tasks, with the one proposed in [Haji-Abolhassani and Clark, 2012b], that uses a TSHMM with no LM (i.e., with a unigram LM), in four different dictionary sizes. We denote the TSHMM that uses the lexicon information by +LM and the TSHMM that disregards any high-level information by -LM. We created four sets of dictionaries of 26, 52, 104 and 312 English words using the Carnegie Mellon pronouncing dictionary (CMPD) [Weide, 2005]. All

bias). The 3-letter words were selected so that there was no repetition of characters in them. At the beginning of every experimental session, we calibrated the eye tracker by having the participant look at a 16-point calibration display that extended to 10×10 degrees of visual angle (the area covered by the calibration grid is stretched beyond the stimuli which spans a 6.6×6.6 degrees of visual angle).

An eye tracker (ISCAN RK-726PCI) was used to record

426 dictionaries were built so that they all include all the 427
428 words of the smaller dictionaries. The words were selected 429
430 randomly from the CMPD and the words length varied 431
432 between three to five characters. The language model was 433
434 also created using the CMU-Cambridge toolkit [Clarkson 435
436 and Rosenfeld, 1997] by extracting language models from 437
438 the words in dictionaries. 439

440 In order to train the TSHMMs, we have to adjust the 441
442 mean vector of the 2-D Gaussians according to the training 443
444 character so that it aligns with the center of character loca- 445
446 tion. According to [Rabiner, 1990] a uniform (or random) 447
448 initial estimation of initial state and transition probabilities 449
450 (II and A) is adequate for giving useful re-estimation of 451
452 these parameters (subject to the stochastic and the non-zero 453
454 value constraints). Thus, we set a random initial values for 455
456 the parameters in the generic HMM and run the Baum- 457
458 Welch algorithm on the training set to obtain the final 459
460 TSHMM [Huang et al., 1990]. We also used a technique 461
462 called parameter tying [Rabiner, 1990] to force a unique 463
464 task and stimuli independent covariance matrix across all 465
466 of the Gaussian distributions in the mixtures. Thus, we 467
468 can build the word model for the test data by dynamically 469
469 changing the means of the states according to the character 470
471 locations of the characters and using the estimated variances 472
473 of characters. 474

475 Figure 7 shows the accuracy of word inference using 476
477 TSHMM with LM (+LM) and TSHMM without LM (-LM) 478
479 methods ranging over four dictionary sizes. As expected, 480
481 the +LM performs better than -LM due to the fusion of 482
483 information provided by the LM. The table below the figure 484
485 shows the accuracy and the standard error of the mean 486
487 (SEM) of the corresponding bars. For each bar we ran a 488
489 10-fold cross validation on our database of 145 trajectories 490
491 in order to define the training and test sets and used the 492
493 same epochs across all the methods. 494

495 ACKNOWLEDGMENT 496

497 The authors would like to thank Fonds de recherche 498
499 du Quebec - Nature et technologies (FQRNT) and Natural 500
501 Sciences and Engineering Research Council of Canada 502
503 (NSERC) for supporting this work. 504

505 REFERENCES 506

507 G.T. Buswell. *How people look at pictures: A study of the* 508
509 *psychology of perception in art*. Chicago: University of 510
511 Chicago Press, 1935. 512
513 M.S. Castelhana, M.L. Mack, and J.M. Henderson. Viewing 514
515 task influences eye movement control during active scene 516
517 perception. *Journal of Vision*, 9(3):6, 2009. 518
519 J.J. Clark and J.K. O'Regan. Word ambiguity and the 520
521 optimal viewing position in reading. *Vision Research*, 39 522
523 (4):843–857, 1998. 524
525 P. Clarkson and R. Rosenfeld. Statistical language modeling 526
527 using the cmu-cambridge toolkit. In *Fifth European* 528
529

Conference on Speech Communication and Technology, 1997.

B. Fischer and H. Weber. Express saccades and visual attention. *Behavioral and Brain Sciences*, 16:553–553, 1993.

GC Gilmore, H. Hersh, A. Caramazza, and J. Griffin. Multidimensional letter similarity derived from recognition errors. *Attention, Perception, & Psychophysics*, 25(5): 425–431, 1979.

A. Haji-Abolhassani and J.J. Clark. Realization of an inverse yarbus process via hidden markov models for visual-task inference. *Journal of Vision*, 11(11):218–218, 2011a.

A. Haji-Abolhassani and J.J. Clark. Visual task inference using hidden markov models. In *Proceedings of the 22nd International Joint Conference on Artificial Intelligence*. IJCAI/AAAI, 2011b.

A. Haji-Abolhassani and J.J. Clark. A computational model for visual-task inference in pop-out and conjunction search. *submitted to Vision Research*, 2012a.

A. Haji-Abolhassani and J.J. Clark. Visual task inference in conjunction search using hidden markov models. *accepted for presentation in Vision Sciences Society 2012 - to appear in Journal of Vision*, 2012b.

H.V. Helmholtz. *Handbuch der Physiologischen Optik*, Dritter Abschnitt, 1896.

J. Hu, M.K. Brown, and W. Turin. HMM based on-line handwriting recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 18(10):1039–1045, 1996.

X. Huang, A. Acero, H.W. Hon, et al. *Spoken language processing*. Prentice Hall PTR New Jersey, 2001.

X.D. Huang, Y. Ariki, and M.A. Jack. *Hidden Markov Models for Speech Recognition*. Edinburgh University Press, 1990.

G. Keren and S. Baggen. Recognition models of alphanumeric characters. *Attention, Perception, & Psychophysics*, 29(3):234–246, 1981.

V. Nair and J.J. Clark. Automated visual surveillance using hidden markov models. In *International Conference on Vision Interface*, pages 88–93, 2002.

L.R. Rabiner. A tutorial on hidden Markov models and selected applications in speech recognition. *Readings in speech recognition*, 53(3):267–296, 1990.

A.M. Treisman and G. Gelade. A feature-integration theory of attention. *Cognitive psychology*, 12(1):97–136, 1980. ISSN 0010-0285.

R. Weide. The carnegie mellon pronouncing dictionary [cmudict. 0.6], 2005.

A.L. Yarbus. Eye movements during perception of complex objects. *Eye movements and vision*, 7:171–196, 1967.

S.J. Young, NH Russell, and JHS Thornton. Token passing: a simple conceptual model for connected speech recognition systems. *Cambridge University Engineering Department*, pages 1–23, 1989.