

CONSTRUCTIVE VISUAL IMAGERY AND PERCEPTION

Arthur M. Farley
University of Oregon
Eugene, OR 97403

Introduction

One active area of artificial intelligence research is the inquiry into the nature of human cognition. One aspect of this investigation is the attempt to embody theories of cognition in the form of programs which simulate both general characteristics and specific instances of observed cognitive behavior. This paper is a report of one such effort in this area of artificial intelligence viewed as theoretical psychology. It is based upon research conducted and reported as a Ph.D. thesis in computer science at Carnegie-Mellon University (1).

Classical or *atomistic* theory proposed that the visual perception of form is an "unconscious conclusion" realized by "unconscious inferences" which are based upon the values of the smallest discernible or homogeneous patches of the stimulus (2). The Gestalt theory of visual form perception developed as a reaction to the failure of the classical approach to adequately account for the effects of context upon the valuation (interpretation) of any part (atom) of the stimulus. Gestalt theory defines several organizational principles (laws of proximity, continuity, symmetry, simplicity) which are applied to the whole perceptual field to produce the visual form perception (3). Gestalt theory, like atomistic theory, has proven to have its weaknesses. The laws have proven difficult to specify in quantitative or operational terms. Perception of partial figure regions can have significant effects upon complete figural perception (4,5). Most notably, Gestalt theory ignores the fact that visual form perception normally involves multiple fixations of the eye and so provides no means for the integration of information from successive differing views.

The constructive theory of visual form perception has developed as an alternative to the classical and Gestalt approaches. This theory proposes that an internal representation of the visual field is constructed by the integration of a succession of views of (fixations upon) the environment. This representation is both guide for and product of visual form perception. Hebb (6) began the modern psychological formulation of this theory, describing "cell assemblies" joined together by (into) "phase sequences" as its basic functional elements. Hochberg (7,8) has recently continued the investigation. He proposes "schematic maps" as the underlying structural organizations which make possible the selective attention to and the successive integration of the visual environment.

The research which is reported here is a further investigation and specification of the constructive theory of visual form perception. More specifically, the goals have been: (1) to investigate the nature of the processes and

memories which are involved in the fixation and integration of successive views of the environment; (2) to investigate the nature of the internal representation (symbolic visual image) which is capable of embodying the necessary partial and complete perceptions; (3) to specify the results of the investigations in the form of an operational, computer-implemented visual imagery and perception system (VIPS). VIPS is such a program which has been implemented in LISP 1.6 on the CUM POP-10.

Motivation for the two investigative goals is abundant. The need is best expressed by the following two statements:

"There is little evidence to guide our thinking on how these integrations and constructions take place. In fact, . . . , little attention has been paid to how such processes occur at all." (9,p174)

"What we need is a set of operations for defining and studying the kind of visual storage that will build up the structures of perceived forms out of momentary glimpses" (7,p322)

The third goal, that of computer implementation of the theory, had two primary motivations. One was to force an operational specification of the theory, something which is often elusive for concisely stated, intuitively understood, descriptive theories. The other was that features of the implementation could be expected to be naturally applicable in extension to the realization of an adequate, generalized computer vision system.

THE EXPERIMENTS

Two experiments were conducted to provide data from which to infer characteristics of the visual image representation and rules of the perceptual processes. The first experiment presented subjects with a task situation which forced them to perform a perceptual activity over an extended (cognitive) time frame. VIPS has been implemented to explain this behavior. The second experiment presented task situations nearer to that of "normal" visual form perception. These results are considered in light of the theory embodied by VIPS and as bases for its necessary extension and modification.

The main task of Experiment I presented the subject (4 subjects) with a line drawing taped to a table which was covered by a large paper mask with a hole in it. The hole (of approx. 3 degrees visual angle) allowed the subject to view at most one vertex at a time. The subject's task was to move the head about the drawing until being capable of verbally describing and drawing the picture. The subject was also instructed to "think aloud"

during the hole movement sequence. This task is an extension and modification of a partial, sequential viewing task described by Hochberg (),

Transcriptions of the video-taped protocols served as a basis for perceptual process rule specification, for the structure of recognition long term memory, and, together with the verbal descriptions and drawings, for image representation specification. Figure 1 is the initial segment of a transcribed protocol. B' illustrates the view through the hole (circle) as then seen by the subject. If this is a vertex, a label number is given in the lower right. C' indicates the direction of hole movement. If any, by arrow and symbol. D' is the transcribed verbalization which occurred at that point. This example illustrates the type of behavior this task situation produced.

	B'	C'	D'
V1			looks like the corner [S1] the part right here looks like the corner of a uh [S2] of a square [S3] now, lets see if I can... [S4]
V2		RT	if it is to be a square [S5] this would be the bottom of the square [S6]
V3			what [S7] looks like we might have a triangular object [S8] this would be
V4		UL	the second side of the triangle [S10]
V5		UL	(slowly)
V6		UL	(through vertex 3)
V7		UL	
V8		DO	and this is the third side [S11]
V9			okay [S12]

Figure 1

Experiment II was conducted upon a corneal reflecting eye movement tracking and recording (video-taping) system. Two subjects participated; both had been subjects for Experiment I. In the first part, a line drawing (of approx. 20 degrees visual angle) was presented for 250 to 350 milliseconds while the subject fixated a pre-set point of the visual field. Immediately upon removal of

the picture, the subject proceeded to verbally describe and to draw what he now knew of the line drawing. This data served as basis for inferring the processing that an Initial fixation receives and characteristics of the resultant representation.

In the second part, the subject was again given a pre-set Initial fixation point. A line drawing (approx. 20 degrees) was presented, and the subject was allowed to view (scan) the line drawing until being capable of verbally describing and drawing it from memory. The picture was removed at the subject's signal and the subject proceeded first to describe, and then to draw, the line drawing. Eye movements were video-taped and the verbal descriptions were audio-taped. This data was transcribed as sequences of fixations specified by location and duration. Figure 2 shows one such transcription. Field A' is fixation number, B' is fixation duration. initial segment of the sequence has been connected by dotted lines. The data has been used to Judge the feasibility of accommodating the theory of VIPS to account for eye movement data.

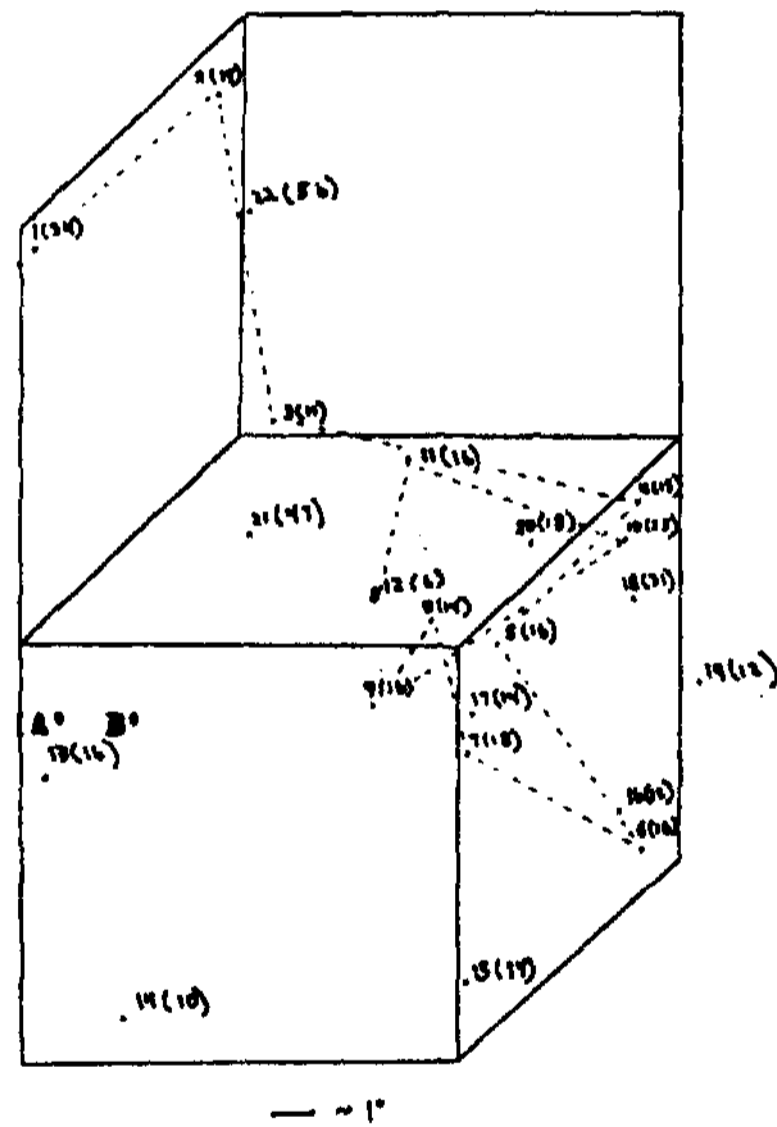


Figure 2

THE IMAGE REPRESENTATION

Imagery has recently returned as a concept under investigation in cognitive psychology. One result of this renewed interest is the need for an adequate scientific definition of imagery. An image is defined here to be an internal, semantic, symbolic representation of information which is capable of determining (guiding) behavior and

which has an Internal modality characteristic. The Internal modality characteristic distinguishes Images from other forms of Internal symbolic representation. It requires that the Image be structured so that It can be (Is accessed by processes isomorphic to those which access the external environment for the specific sensory modality (10). The units of representation, the symbols and relations, are also modality specific, for vision this being visual feature symbols and spatial relations. By further specification, a visual image is an Internal symbolic representation of visual feature and spatial relation Information which Is structured so as to be straightforwardly retrievable by Internal processes Isomorphic to those which access the visual environment during visual perception. A visual Image Is capable of guiding motor and cognitive behavior with regard to the visual environment that It represents. A primary example of such behavior is visual form perception. On the other hand, visual form perception Is a primary source of visual images. Thus, the perceptual visual Image serves as guide for its own construction.

The basic unit of meaningful visual Information within VIPS is the image chunk. An image chunk is a semantic structure of Interrelated symbolic elements, which represent a visual concept. image chunks are the basic structural units from which more complex Images are then constructed. VIPS uses five types of image chunks (concept types) to construct Its representations of the line-drawing environment, these being the VERTEX, OBJECT, LINE, SIDE, and FACE types.

An image chunk consists of one Chunk Header element, one or more Position elements, and several image Body elements. Properties of the Chunk Header serve to Indicate the chunk's type and general characteristics of the current instance. The Current Reference property always references one image Body element, thereby affording a means of access to the chunk's representational structure. A Position element serves to Indicate the perceived location of the Image Body elements which reference It. This position Is In terms of a seven-by-seven coordinate grid of locational areas. All Image Body elements reference a Position element, binding the Image Body structure, and thus the visual Image, to locations in perceptual (imaginal) space. Perceptual space is not retinally-based, but rather Is bound to the area of current Interest within the visual field. This Indirectness plays a significant role In the maintenance of a stable visual world In spite of the differing retinal states which occur due to fixation changes during visual form perception.

A chunk's image Body elements form a non-hierarchical symbol structure which embodies the spatial configuration of visual features of the visual concept represented by the chunk. This structure is a doubly linked circular list. The relational links between Image Body elements have directional or space-traversal (direction and distance) meaning. Thus, the allowed means of accessing the Image, which Is by traversal of image Body elements according to existing structural links, is isomorphic In meaning to a visual search, or scan of the external visual environment. This feature of the Image representation

is the embodiment of the Internal modality characteristic for visual Imagery.

Five classes of image Body elements (XIT, ANGLE, INTERNAL, END, and QUICKSEE) are used to represent the concepts embodied by Image chunks. Each element class has an associated set of features and relations that it can embody which in turn defines possible roles for It In the Image Body structure. Figure 3 Illustrates the use of these element classes in representative examples of the different chunk type Image Body structures. A VERTEX chunk consists of alternating XIT and ANGLE elements. An OBJECT chunk consists of "corner configurations" made up of INTERNAL and ANGLE elements. An INTERNAL differs from an XIT element In that It has only one Internal vertex direction link to an ANGLE element. This Image representation feature embodies the figure-ground phenomenon of visual form perception and imagery. To traverse around the "outside" of an object corner, another related chunk must be accessed (attended). The LINE and SIDE chunks introduce the use of the END element. A SIDE chunk Is always associated with an OBJECT chunk.

A visual Image Is a structure of Interrelated Image chunks. An Image Body element of one Image chunk may be linked to one of another chunk by one of two basic relation types. One type has space traversal meaning, relating the two Image Body elements in terms of a direction and distance. The other type Is an equivalency relation. This type of link relates two Image Body elements of different chunks which represent the same aspect of the external visual field. For example, a line segment shared by two adjacent objects Is redundantly represented by both OBJECT chunks. INTERNAL elements embodying the shared line segment "side" of the objects are linked by inter-chunk equivalency relations. (See bottom Figure 3) Note that the line segment Is redundantly represented In both OBJECT chunks.

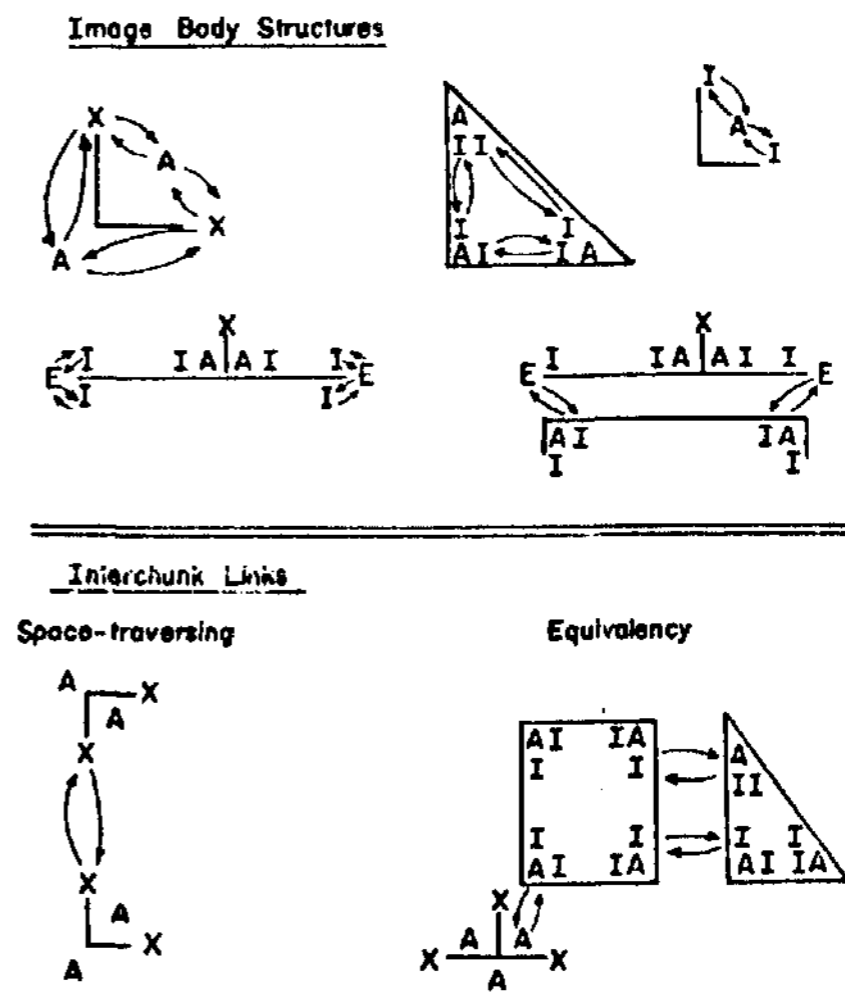


Figure 3

The symbolic image representations of VIPS differ from the more semantic representations of visual information proposed by Winston (11). His representations do not satisfy the image definition given above. Such relations as "SUPPORTED-BY" and those accompanied by satellites, such as "MUST-NOT-ABUT", *are* not spatial or equivalency relations. Such semantic relations do appear to play a role in perception, possibly in the incorporation process which forms more complex, "deeper" structures from interrelated image chunks. The process could provide perceptual hypotheses based upon the long term semantic knowledge to help realize the goal of efficient representation. Also, VIPS does not currently employ a "SCENE" element, but it would exist in the guise of a contextual element (1).

THE PERCEPTUAL SYSTEM

The perceptual system of VIPS consists of six memories and four processes, as shown by Figure 4. In the figure, an arrow from a process to a memory indicates that the process can alter the memory's contents, while an arrow from a memory to a process indicates that the process can access the memory's contents. Characteristics of the system at the architectural level reflect relevant data and theoretical proposals of cognitive psychology.

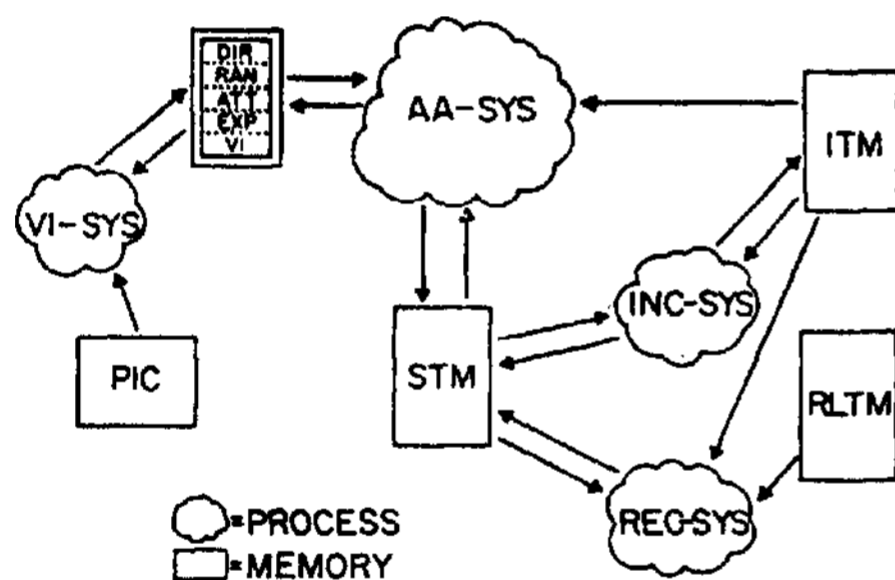


Figure 4

The Memories

Since perceptual activity involves the interaction of organism and environment, VIPS must represent both. PIC (Picture) is the environment, defined in terms of the visual modality. As such, it consists of a list of vertex feature lists. Each vertex feature list represents the inferred detectable visual features at one of the line drawing vertices (the observed points of "hole fixation"). These vertex feature lists are interrelated by links of directional and distance meaning. The Current Picture Pointer (CPP) is associated with PIC and references one vertex feature list, representing the current hole position.

The five cells of the Visual Register (VR) serve as communication registers between the

VI-SYS and AA-SYS processes. The cells of VR are set prior to VI-SYS activation by AA-SYS. VI-SYS accesses PIC according to the VR cell specifications and accordingly alters the cells of VR prior to returning control to AA-SYS by deactivation.

The VI (Visual Information) cell is the iconic visual image. It is constructed automatically as a result of the new hole positioning. The representation of visual information in VI is that of an image chunk, though a simple feature list appears favorable. It is as yet unrelated to any existing image contents produced by prior perceptual activity. Its contents and structure are not affected by any active perceptual goal.

When specified, EXP (EXPection) and ATT (ATTention) make possible the application of pre-attentive functions by VI-SYS to the newly accessed PIC information. EXP can be specified either as an angle code (when ATT is also specified) or as a vertex type and specification. Upon VI-SYS deactivation, EXP will be YES, NO, or CON (CONTAINED), indicating the relationship of expectation to realization. ATT can be specified as a direction symbol. If specified it enables VI-SYS to move through any encountered vertex which is straight (has 180 degree angle) on the ATT side of the line being traversed.

DIR (DIRection) specifies the direction of the hole movement to be effected by VI-SYS, RAN (RANge) specifies the range of that move. RAN is a value returned by VI-SYS in VIPS (hole movement); during saccadic (ballistic) eye movement behavior, it may be pre-specified by AA-SYS. The values of these cells are incorporated into the image by AA-SYS. This reflects the intrinsic role of motor (efferent) activity and symbols in visual perception and imagery.

Short Term Memory (STM) consists of an ordered list of nine chunks, being the limited amount of active memory available to the perceptual (cognitive) system for image construction (12). An STM chunk is an image chunk or an image chunk with a special type element (GOL, LAST, COM, or OGOL) appended. STM chunks are accessed according to image body element containment of special element appendment. STM is the memory in which the visual image (perception) is constructed. STM and VR *are* the active memory components or "mind's eye" during visual form perception. During recall and drawing involving imagery, only STM is active as the "mind's eye".

There are three components of long-term memory in VIPS. LTM is the universe of symbols used in VIPS. The symbols are interrelated, forming a semantic network of symbols and relations. The use of any symbol in an active memory or by a process rule is an activation of that LTM element. (In the implementation in LISP 1.6, it is a pointer to that element.) As such, any active memory entity (VR cell or STM chunk) is only pseudo self-contained. A symbol instance implies the possible use by the active process of all the symbol's LTM relations and related symbols. A symbol remains effectively imbedded in the LTM structure when activated.

RLTM (Recognition Long Term Memory) is a modified n-any discrimination net. The discrimination structure is determined by angle symbols and the objects which angle configuration determine. The discrimination net is not used to directly guide perceptual activity. Rather, it is used to name completed object images (if possible) or to activate an available image generating process to complete a hypothesized object image when adequate partial information is available. RLTM is accessible only by REC-SYS and is therefore only traversed during REC-SYS activation.

ITM (Intermediate Term Memory) is the memory into which selected STM chunks are incorporated. Meaningful confirmed results of the perceptual activity are transferred to this memory, as the current meaningful processing context. Its contents are recallable into STM for use by the perceptual process. This recall requires a full system cycle thus the contents are not immediately accessible. As a chunk is incorporated into ITM new relational properties are added to its chunk header element. These indicate the chunk's immediate image context, being the image chunks to which its image body elements are linked. Temporal relations between OBJECT chunks are added also. Thus, the final perceptual image, which is found in ITM, is a heterarchical symbol structure.

The Processes

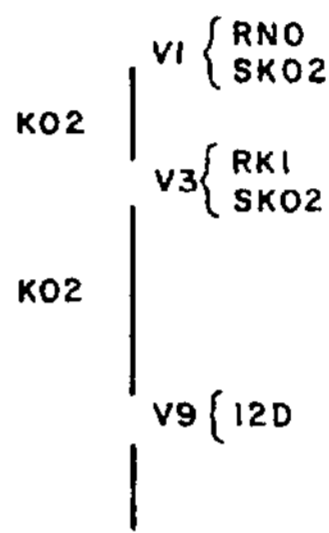
All four processes are implemented in the form of production systems (13). A production system consists of an ordered list of condition-action pairs. The system is cyclic in operation. With each cycle, the first rule of the active process which has its condition-half satisfied is said to "fire", resulting in its action-half being executed. Rule firings depend upon and can alter only active memory contents. This reflects the contextual nature of human cognition.

The Assimilation-Accommodation process (AA-SYS) is named for the two basic classes of behavior associated with it (14). Upon VI-SYS deactivation, this process either can assimilate the new contents of VR into the image in STM or must accommodate that current image in light of the conflicting contents of VR. AA-SYS is the "main" process. It is the source of most goals. It sets VR and activates VI-SYS so to access new PIC contents. It activates REC-SYS to aid in its image construction activity. It activates INC-SYS to incorporate satisfactorily confirmed, meaningful segments (chunks) of the image into ITM.

As the primary source of current goals and their transitions, AA-SYS embodies perceptual strategy in the VIPS implementation. Two overall strategies have been inferred from the protocol data, resulting in two corresponding AA-SYS implementations. One strategy is that of successively recognizing objects until having represented the whole line drawing. The other is that of first attempting to scan and represent the entire line drawing outline and then linking up unknown inward directed exits. The second strategy diverts to the first in a number of specified circumstances.

Each rule of AA-SYS, REC-SYS, and INC-SYS has the name of the current goal as the primary

consideration of its condition-half. Visual form perception is a goal directed activity which consists of a sequence of goal-related episodes. Verbalizations and repeated patterns of hole movements were indications of this that can be seen in the protocols. The current goal is a primary determinant of what image chunks are attended, what external information is to be accessed (fixated) next, and into what class of image structures the new information will be integrated and by which it will then be represented. Figure 5 is a goal episode chart of the goal sequence as inferred for the protocol segment of Figure 1. Goals to the right (in brackets) of a protocol frame number are active at (during) that "hole fixation". A goal to the left of a vertical line is active during the hole movement sequence occurring between the protocol frames indicated. The resultant perception is the product of the interaction (or interference) between cognitive goals, partial perceptions, and visual environment.



RNO—Recognize New Object
SKO2—Start Known Object 2-Dimensions
KO2—Known Object 2-Dimensions
RKI—Rerecognize, Known Interrupted
I2D—Incorporate 2-Dimensional

Figure 5

The output produced by the VIPS production systems is a trace of rule firings and resultant active memory contents. Figure 6 presents the VIPS trace corresponding to the human behavior transcribed as frames VI and V2 of Figure 1. VI has been set to the initial vertex and STM is empty. Rule RAO of AA-SYS fires, assimilating the icon into STM as a new image chunk. The goal RNO is generated by the rule, and a new chunk is created, marked by the special element GOL. This new chunk is to be the sight of object image construction during the ensuing REC-SYS activation. Rule RECO of REC-SYS next fires and traverses the vertex chunk in STM (marked by LAST), beginning construction of the new object accordingly. RECI next fires as there is no further relevant image information in STM. The UOX property of the node currently referenced in RLTM activates a process to complete a square-or-rectangle image, as is suggested by the right angle in STM. REC-SYS then deactivates, control returning to AA-SYS. Rule SKO2I fires, using the new object image to specify the contents of the visual register before activating VI-SYS. Rule VI9 fires, moving the CCP "right" to the feature list for vertex 3.

The rule constructs the icon, provides range information, and returns NO in EXP.

```

VI IS (V3 P4 X5 A8 X7 A6)

GOAL-IS: NIL
2,0 STM: (0 0 0 0 0 0 0 0)

ACT #1: AA-SYS.

RULE RAAD: BEF FIRST VISUAL INFO.

GOAL-IS: RNO
3,0 STM: (GOL )
(LAST V3(CR X7 PNUM P4 VSPEC (RTA DL) SPEC V2 TYPE VERTEX)-
(P4 X5 A8 X7 A6) )
(0 0 0 0 0 0)

ACT #9: REC-SYS.

RULE RECO: BEGIN RECOGNITION OF NEW OBJECT.

GOAL-IS: RNO
4,0 STM: (GOL 010(DIM TWO NUMS ONE TYPE OBJECTXP4 I11 A12 I-
13) )
(COM V3(CR X5 PNUM P4 VSPEC (RTA DL) SPEC V2 TYPE VERTEX)-
P4 X5 A8 X7 A6) )
(0 0 0 0 0 0)

RULE REC1: NO MORE MEMORY INFORMATION.

GOAL-IS: SKO2
5,0 STM: (GOL 010(NAME SQUARE-OR-RECTANGLE DIM TWO TYPE OBJE-
CTXP4 I11 A12 I13 I14 A15 I16 I17 A18 I19 I20 A21 I22) )
(V3(CR X5 PNUM P4 VSPEC (RTA DL) SPEC V2 TYPE VERTEXXP4 X5 A8-
X7 A6) )
(0 0 0 0 0 0)

ACT #11: AA-SYS

RULE SKO21: START LOOKING FOR KNOWN OBJECT INFO.

ACT #23: VI-SYS.

RULE V19: EXP & ATT & NO STRAIGHT SIDES FOUND

DIR IS RT
RAN IS NIL
EXP IS RTA
ATT IS UP
VI IS NIL
AT VERTEX VER3
DIR IS RT
RAN IS LO
EXP IS NO
ATT IS UP
VI IS (V24 P25 X26 A31 X30 A29 X28 A27)

```

Figure 6

Figure 7 presents pictorial representations of the image chunks created by VIPS in its account of the remainder of the protocol section of Figure 1. The NO in EXP causes object image accommodation. The new vertex is assimilated as a chunk, the square image is altered to include only the right and acute angles, REC-SYS is reactivated with goal RNI and a triangle image is completed, as suggested by the partial object image in STM. As ATT will be specified as the direction to the inside of the triangle, VI-SYS "moves" through vertex 8 during its ensuing activation. Since the new vertex is not simple, it is assimilated into STM and linked by equivalency to the triangle, a SIDE chunk is created with a QUICKSEE element included, which notes only the number of exits in the by-passed vertex. Vertex is next re-encountered in the system's attempt to confirm the triangle image. The vertex agrees with that expected so it is not assimilated into STM, but is "seen" as part of the now completely confirmed triangle image. Being completely confirmed, the

triangle image and its immediate image environment are incorporated into ITM. In general, VIPS accounts for over 80 percent of the hole movements. The system consistently has sufficient image chunks in STM to be a source of the observed verbalizations.

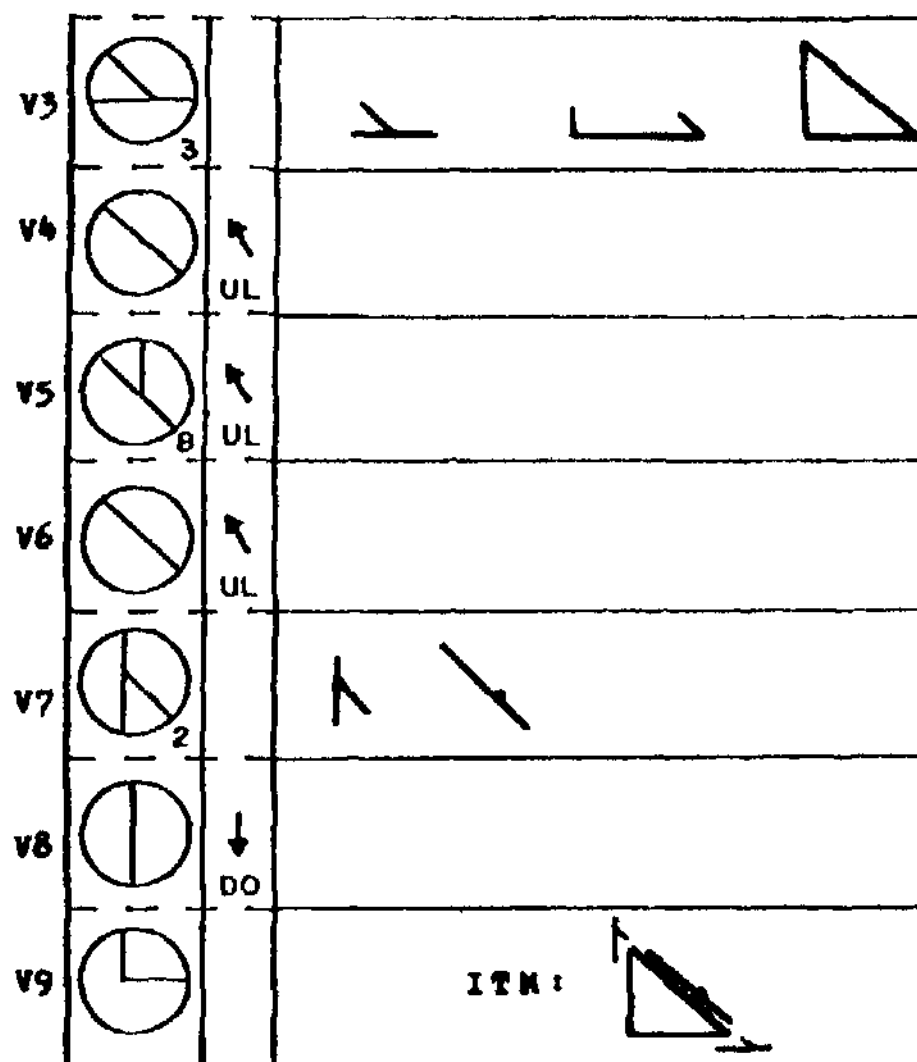


Figure 7

THE INITIAL FIXATION AND EYE MOVEMENTS

The behavior which was observed and recorded as data in Experiment II has been considered only at a general, descriptive level. No perceptual system has yet been implemented which produces corresponding behavior and which could then serve as an operative theoretical explanation of the behavior. This general consideration of the data has been favorable to the proposition that VIPS can be transformed into a sufficient theory of that behavior.

The verbal descriptions and partial drawings obtained as data in Part I of Experiment II provide the basis for inferring the state of the "mind's eye" following processing of the initial fixation by the perceptual system. The data indicate that the initial fixation serves the perceptual system as the source (basis) of inferences and hypotheses. General characteristics of the total extent of the line drawing are consistently inferred, thus determining (delimiting) the current area of perceptual interest. Object hypotheses are generated for regions (areas bounded by spatially disjoint features) within that extent of the visual field. Such processing of the initial fixation is consistent with an active, assimilation-accommodation theory of visual form perception as is embodied by VIPS.

Much of the Initial processing is based upon visual information lying in the periphery of the visual field. This input must be realized by processes of peripheral vision. The data indicate that peripheral information at eccentricities of up to twenty degrees of visual angle contributed to the Initial Inferences and hypotheses. Some object hypotheses were based upon visual information which must have been entirely peripheral. Though a sufficient basis for object hypotheses, peripheral vision consistently failed to allow subjects to correctly infer the particulars of object interactions.

The hole movement task masked peripheral vision information. Therefore, the representation and utilization of peripheral vision's input are not part of VIPS and are most basic and necessary extensions to the system. A class of elements which embody the position of irregularities or discontinuities in the periphery shows promise as a possible symbolic representation. These elements would occur in VI in addition to the foveally produced feature symbols. The position information would necessarily be retinally-based (relative in direction and distance to the current fixation point). Local differencing operators have been shown to be effective as a means of detecting and locating discontinuities (5). The operators are psychologically and physiologically feasible and readily implementable. These elements appear to be an adequate supplement to foveal (feature) elements as bases for the observed inferences, hypotheses, and eye movement behavior.

The recorded eye movement sequences obtained in Experiment II have been considered in terms of fixation durations and locations. Fixation durations varied from approximately one-eighth of a second to over a second. The variation in fixation length favors an active, assimilation-accommodation theory of visual form perception which predicts differing amounts of information processing per fixation. The initial fixation was significantly longer in duration than those immediately following it, thus upholding its role as a highly processed source of inferences and hypotheses. In general, the eye movement sequences consisted of an initial, object-related, overall scan of the line drawing, perceptual construction relying heavily upon peripheral input, followed by a rescanning which concentrated primarily upon object interaction areas.

For a specific example, consider the eye movement protocol of Figure 2 and Figure 8, which is the resultant verbal description provided by the subject. Part "a" of the verbal description appears to be a conclusion resulting from the initial segment of the protocol (that which is connected by dots). The subject's attention is primarily focused upon the ambiguous (in depth) areas of the line drawing. The remaining sequence of eye fixations indicate that the subject scans the bottom cube and finishes with a long fixation between the remaining faces. One can infer from this sequence the construction of an image sufficient to serve as basis for part "b" of the given description.

that was one of those screwy objects again
 (a) it was a
 it was attempting to be two cubes on top
 of one another
 In order to be able to draw it I said
 there was a normal cube on the bottom
 where you can see the front
 (b) front right hand side and top
 there were two like
 sides
 two like walls
 then there was a wall sitting on top of
 the cubes

Figure 8

TOWARD COMPUTER VISION

VIPS is meant to be an operational theory of human visual form perception, but it has a natural relevance to machine vision research. Three general characteristics of VIPS are applicable, and indispensable I believe, to the eventual realization of an adequate computer vision system. One is the system architecture of VIPS as illustrated in Figure 3. It provides a basis for the appropriate modularization of memory and process involved in perceptual activity. Memories are required (1) to hold the results of the direct analysis of environmental input, (2) in which to construct the perceptual image, (3) to hold confirmed partial and the eventual complete perceptions, (4) to hold applicable recognition (and general) knowledge. Processes are needed (1) to analyze environmental input (i.e., perceptor and gradient operators), (2) to construct perceptual images and select environmental areas for analytic attention, (3) to incorporate confirmed partial perceptions into an overall representation.

The second characteristic is the sequential, localized, in-depth processing of environmental information and the cyclic assimilation-accommodation of the symbolic representation being constructed in light of that processing's results. Early attempts at computer vision found that the information load associated with the storing and processing of a whole scene at one time was astronomical. Since then a main theme of vision research has necessarily been to reduce internal information flow. What VIPS and other eye movement research indicates is that the human vision system faces the same problem and that it approaches (not always solves) the problem by the selective, in-depth processing (fixation) of local areas within the visual environment. The non-fixated region of the visual environment is temporarily given only minor attention and sparse representation. Perception is then accomplished by the sequential application of the in-depth process. With each new fixation, new information is incorporated into the perception under construction. The location of the next fixation is dependent upon the perceptual state following this incorporation. The perceptual state consists of the partial perception in active memory and the current goal.

The third applicable characteristic of VIPS is the goal determined nature of its activity. The current goal plays a role in limiting information flow by influencing what is fixated, how it is represented, and the determination as to when the perceptual construction is sufficient.

thus concluding perceptual activity. Data from visual search tasks illustrate the effects of goals upon the location of fixations during stimulus scanning (16,17). The now classic data presented by Yarbus (18) Indicate clearly the effects that question-answering (Informational) goals have upon the fixation of a presented visual scene. Other research in cognitive psychology has consistently indicated the selectivity of Information representation based upon task-related goals (19).

What this should tell researchers in computer vision is that determining one perceptual strategy and one representational scheme is not the proper target. An adequate vision system will require the availability of several strategy/representation sub-system. Each can then be properly applied in the different goal-related situations which the associated "seeing organism" will find itself. There is *never* "correct perception", only "sufficient perception" for the task at hand.

CONCLUSION

VIPS has been implemented to directly account for four protocols. An evaluation of protocol-program trace correspondence is favorable to the perceptual theory embodied in the system. The element classes of the image representation are shown to form a basis for image generalization. Extensions are discussed which transform the system into a more inclusive perceptual theory. A comparison of studies yields the proposal that human information representation in active memory is flexible, being as suited to the task as possible in light of the existing symbols, relations, and structure of an evolving long term memory upon which all cognition is based.

The research which has been reported here has resulted in a speculative, theoretical specification of the constructive theory of visual form perception. It is not meant to be a conclusion and so this report likewise does not really have one. As an initial implementation, VIPS can now have modifications and extensions applied (some are discussed above) so to better embody the perceptual activity it is meant to explain. Being a theoretical statement, it can serve as a source of experimental questions, the resulting experimental findings being the basis for subsequent specification improvement. Finally, the model also embodies characteristics relevant to the development of a generalized machine vision system.

ABSTRACT

A computer implementation of the constructive theory of visual form perception and imagery is described. An illustrative example of the system in action and a discussion of its correspondence to human behavior are presented. Extensions to the implementation are discussed. Aspects of the model are presented as being relevant to machine vision research.

REFERENCES

- (1) Farley, A. M. VIPS: A Visual Imagery and Perception System; the result of a protocol analysis, Ph.D. Thesis, Computer Science Dept., Carnegie-Mellon University, 1974.
- (2) Helmholtz, H., Handbook of Physiological Optics, Dover reprint, 1963.
- (3) Koffka, K. Principles of Gestalt psychology, New York: Harcourt-Brace, 1935.
- (4) Simon, H. A. An information-processing explanation of some perceptual phenomena. British Journal of Psychology, 1967, 58, 1-12.
- (5) Gregory, R. L. The Intelligent eye. New York: McGraw-Hill, 1970.
- (6) Hebb, D. O. The organization of behavior. New York: Wiley, 1949.
- (7) Hochberg, J. E. In the mind's eye. In Haber, R. N. (ed.), Contemporary theory and research in visual perception, New York: Holt, Rinehart & Winston, 1968, p 309-331.
- (8) Hochberg, J. E. Attention, organization, and consciousness. In Mostofsky (ed.), Attention: contemporary theory and analysis, New York: Appleton, 1970.
- (9) Haber, R. N. & Hershenson, M. The psychology of visual perception. New York: Holt, Rinehart & Winston, 1973.
- (10) Simon, H. A. What is visual imagery? an information processing interpretation. In Gregg, L. W. (ed.), Cognition in learning and memory, New York: McGraw-Hill, 1972, p 183-204.
- (11) Winston, P. H., Learning structural descriptions from examples, MAC-TR-76, Cambridge, Mass: MIT.
- (12) Miller, G. A. The psychology of communication. New York: Basic Books, 1967.
- (13) Newell, A. & Simon, H. A., Human Problem Solving, Englewood Cliffs, N. J.: Prentice Hall, 1972.
- (14) Piaget, J. Structuralism. New York: Basic Books, 1970.
- (15) Uhr, L. Pattern Recognition, learning, and thought. Englewood Cliffs, N. J.: Prentice Hall, 1973.
- (16) Williams, L. G. The effects of target specification on objects fixated during visual search. Acta Psychologica, 1967, 27, 355-360.
- (17) Gould and Dill Eyemovement parameters and pattern discrimination. Perception and Psychophysics, 1969, 6, 311-320.
- (18) Yarbus, A. L. Eye-movements and vision, New York: Plenum Press, 1967.
- (19) Trabasso T. & Bower, G., Attention in learning: theory and research, New York: Wiley, 1968.