

بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ



هوش مصنوعی

فصل ۲۵

ادراک: بینایی ماشینی

Perception: Machine Vision

کاظم فولادی قلعه
دانشکده مهندسی، دانشکدگان فارابی
دانشگاه تهران

<http://courses.fouladi.ir/ai>

هوش مصنوعی

ادراک: بینایی ماشینی



مقدمه

ادراک و عامل

ادراک از حسگرها شروع می‌شود.

حسگر: هر چیزی که بتواند جنبه‌هایی از محیط را ثبت کرده و به صورت یک ورودی به برنامه‌ی عامل ارسال کند.

مثال: حسگر یک‌بیتی ساده: روشن/خاموش، شبکه‌ی چشم انسان

بینایی: مفیدترین حس در برخورد با دنیای فیزیکی

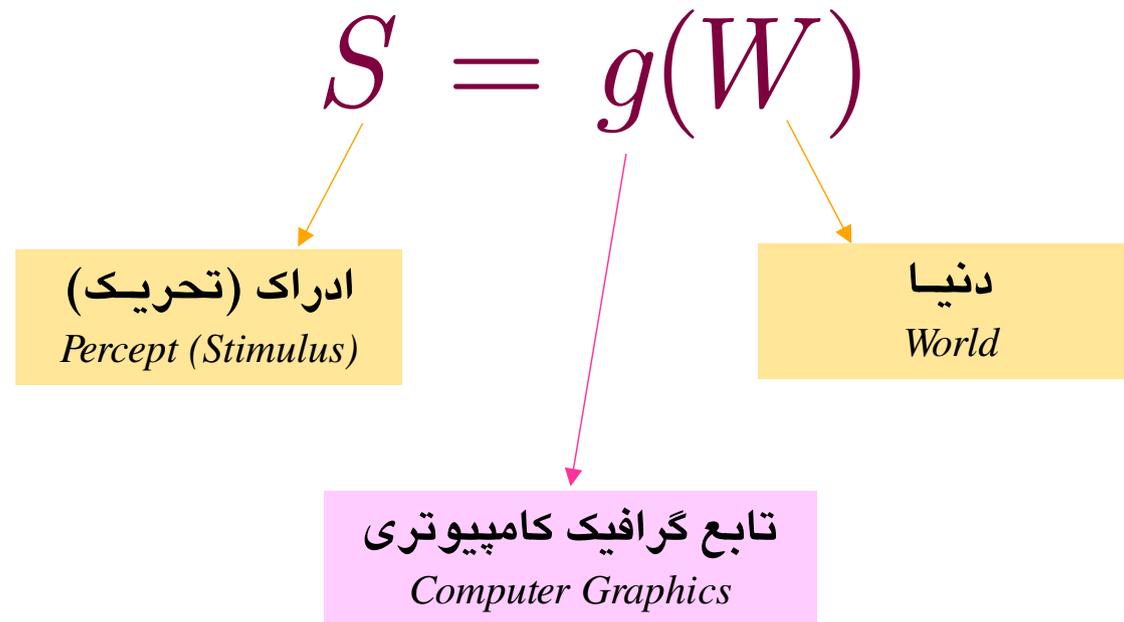
روی‌کرده‌های استفاده از ادراک‌ها توسط عامل

بازسازی <i>Reconstruction</i>	بازشناسی <i>Recognition</i>	استخراج ویژگی <i>Feature Extraction</i>
عامل از روی یک تصویر یا مجموعه‌ای از تصاویر، یک مدل هندسی از دنیا می‌سازد.	عامل بر اساس اطلاعات دیداری و غیره میان اشیایی که با آنها مواجه می‌شود، تمایز قایل می‌شود. بازشناسی می‌تواند تصویر را برچسب‌گذاری کند.	عامل تعداد اندکی از ویژگی‌ها را از ورودی‌های حسی خود شناسایی کرده و آن‌ها را مستقیماً به برنامه‌ی عامل خود می‌فرستد. با ویژگی‌ها به صورت واکنشی برخورد می‌شود.

ادراک: بینایی

گرافیک کامپیوتری

از محرک حسی برای بازسازی دنیا استفاده می‌کنیم:

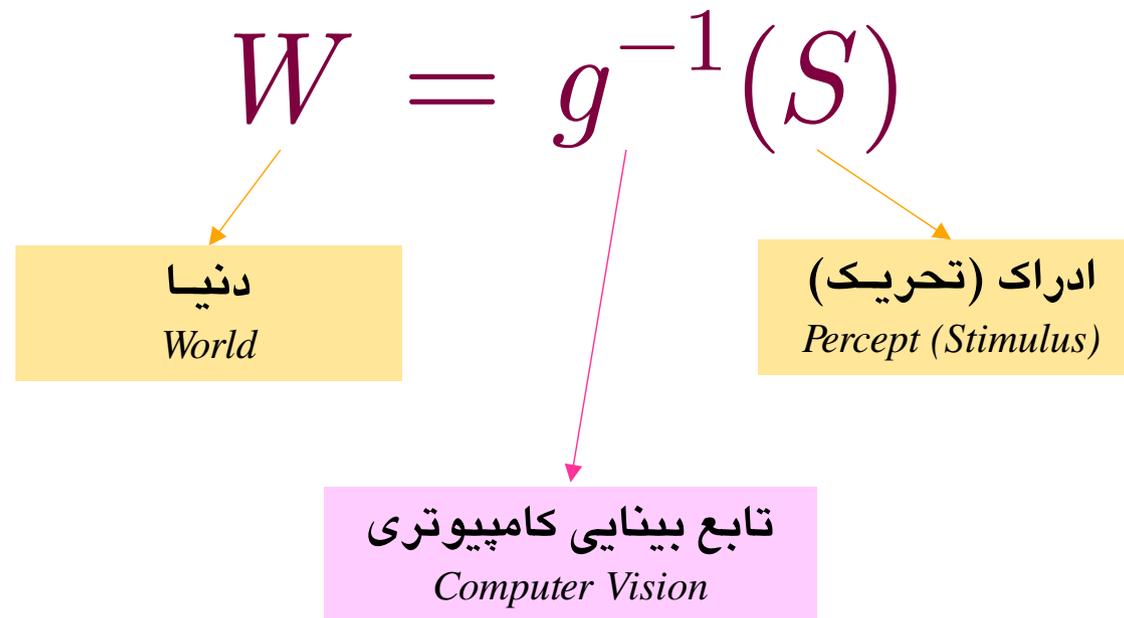


وضعیت دنیا (واقعی یا تخیلی) W را به محرک تولید شده از دنیا S نگاشت می‌دهد.

ادراک: بینایی

بینایی وارون گرافیک است

از محرک حسی برای بازسازی دنیا استفاده می‌کنیم:

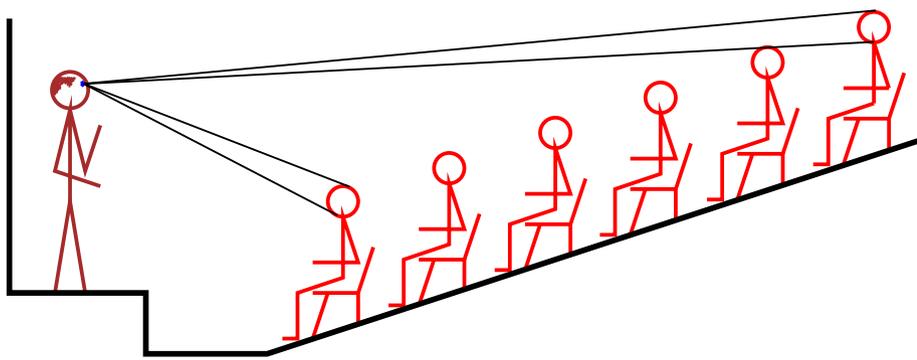
محرک تولید شده از دنیا S را به وضعیت دنیا W نگاشت می‌دهد.اما متأسفانه g وارون مناسبی ندارد!

ادراک: بینایی

مثال (۱ از ۳)

آیا می‌توانیم به بینایی به عنوان وارون گرافیک نگاه کنیم؟

$$W = g^{-1}(S)$$

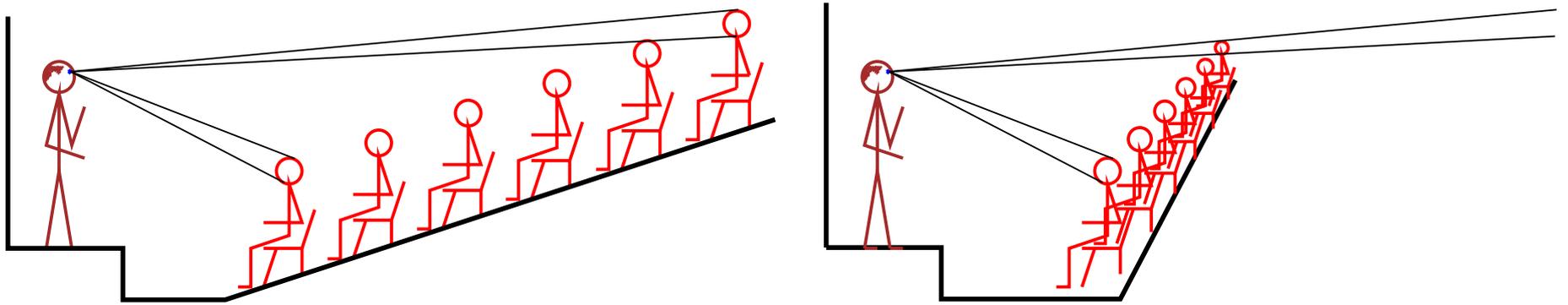


ادراک: بینایی

مثال (۲ از ۳)

آیا می‌توانیم به بینایی به عنوان وارون گرافیک نگاه کنیم؟

$$W = g^{-1}(S)$$



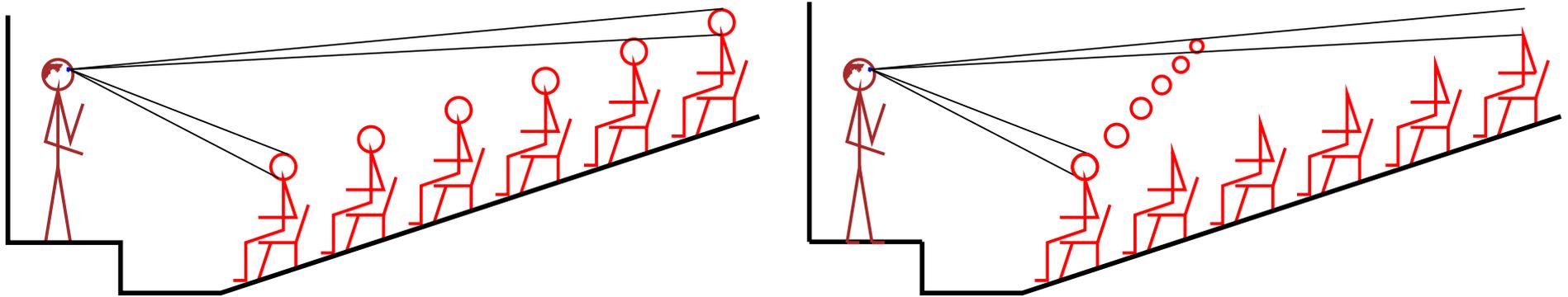
مشکل: ابهام انبوه!

ادراک: بینایی

مثال (۳ از ۳)

آیا می‌توانیم به بینایی به عنوان وارون گرافیک نگاه کنیم؟

$$W = g^{-1}(S)$$



مشکل: ابهام انبوه!

ادراک: بینایی

رویکرد بهتر: رویکرد بیزی

رویکرد بیزی برای استنتاج پیکربندی دنیا

$$P(W | S) = \alpha P(S | W)P(W)$$

گرافیک
Graphics

دانایی پیشین
Prior Knowledge

رویکرد بازهم بهتر: نیازی نداریم که صحنه‌ی دقیق را بازیابی کنیم
فقط اطلاعاتی که برای موارد زیر نیاز داریم را استخراج می‌کنیم:

شناسایی
Identification

بازشناسی
Recognition

دستکاری
Manipulation

ناوبری
Navigation

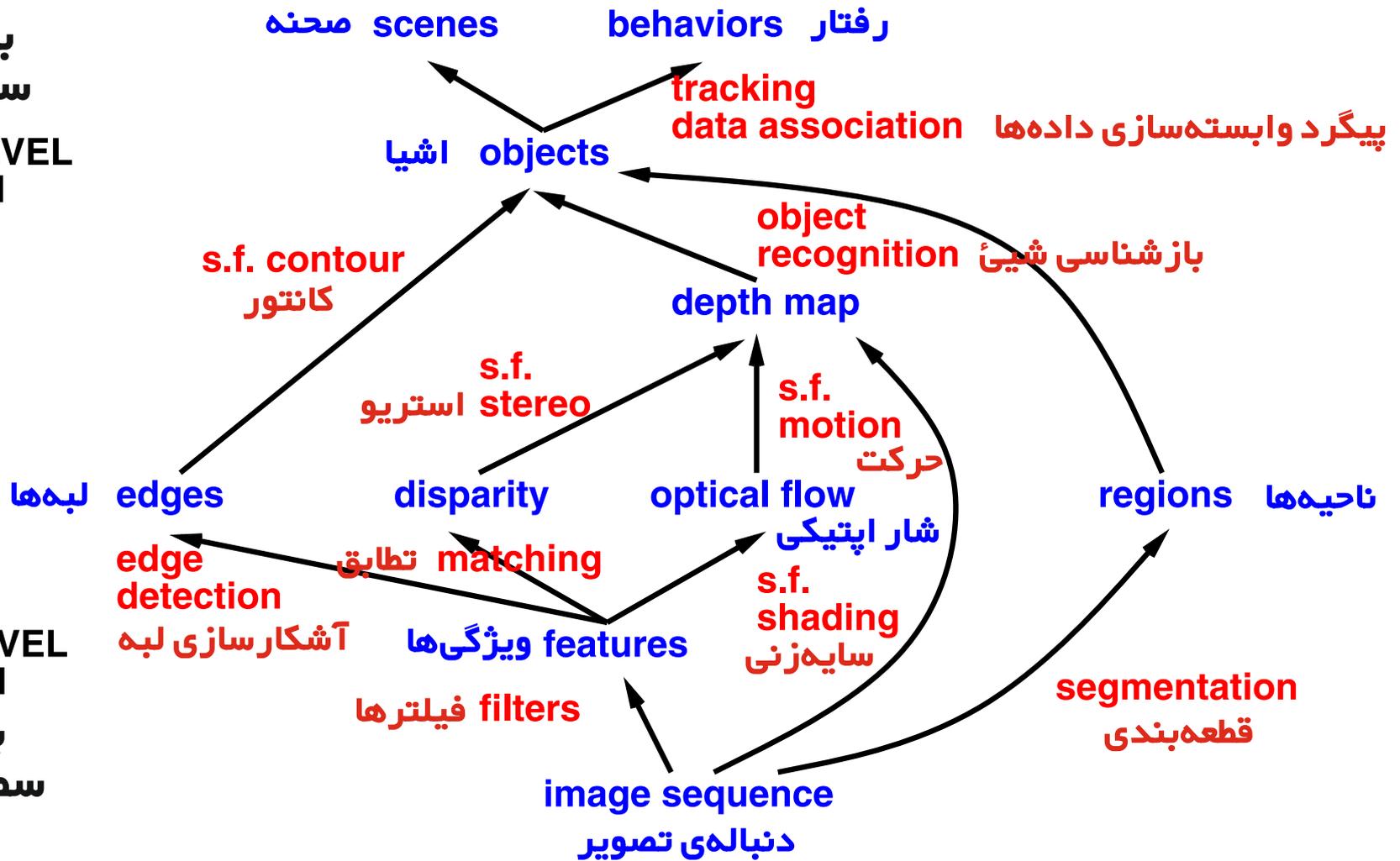
بینایی کامپیوتری

زیرسیستمها

COMPUTER VISION: SUBSYSTEMS

بینایی
سطح بالا
HIGH-LEVEL
VISION

LOW-LEVEL
VISION
بینایی
سطح پایین



بینایی نیاز دارد **اشاره‌های** مختلف را با هم ترکیب کند.

هوش مصنوعی

ادراک: بینایی ماشینی

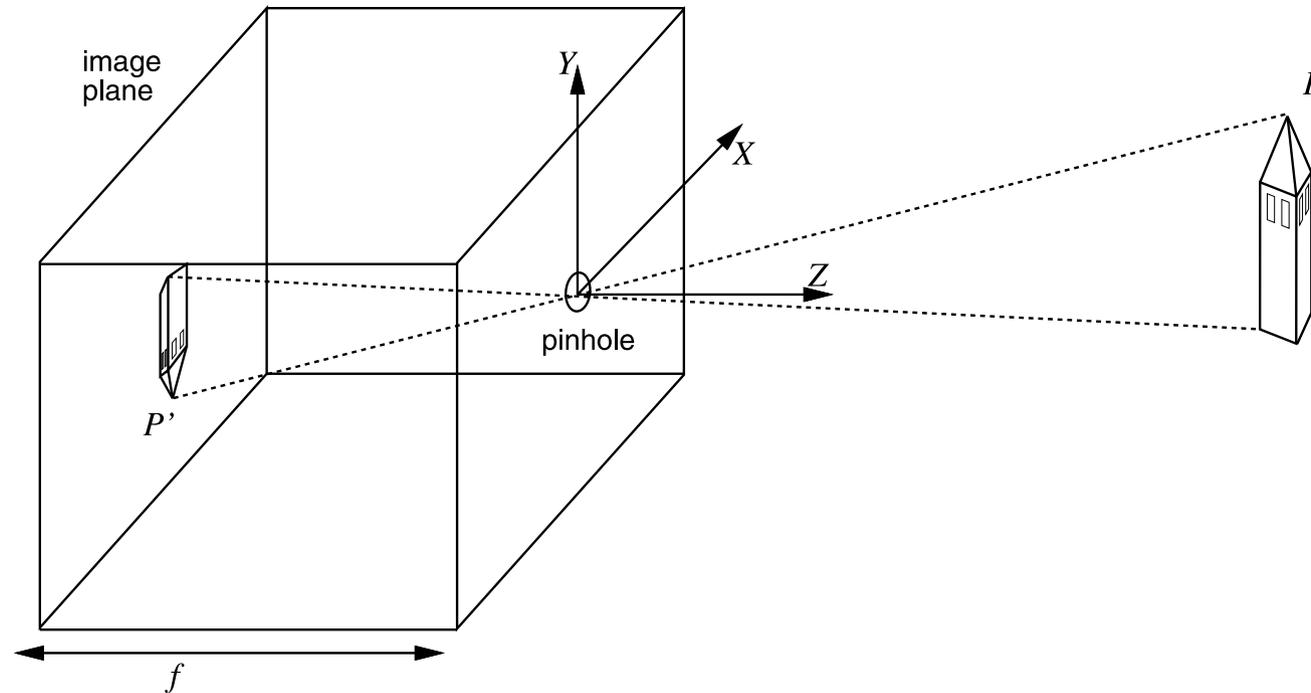
۱

تشکیل
تصویر

تشکیل تصویر

دوربین «سوراخ سوزنی»

IMAGE FORMATION



P یک نقطه در صحنه با مختصات (X, Y, Z) است.
 P' تصویر آن در صفحه‌ی تصویر با مختصات (x, y, z) است.

$$x = \frac{-fX}{Z}, \quad y = \frac{-fY}{Z}$$

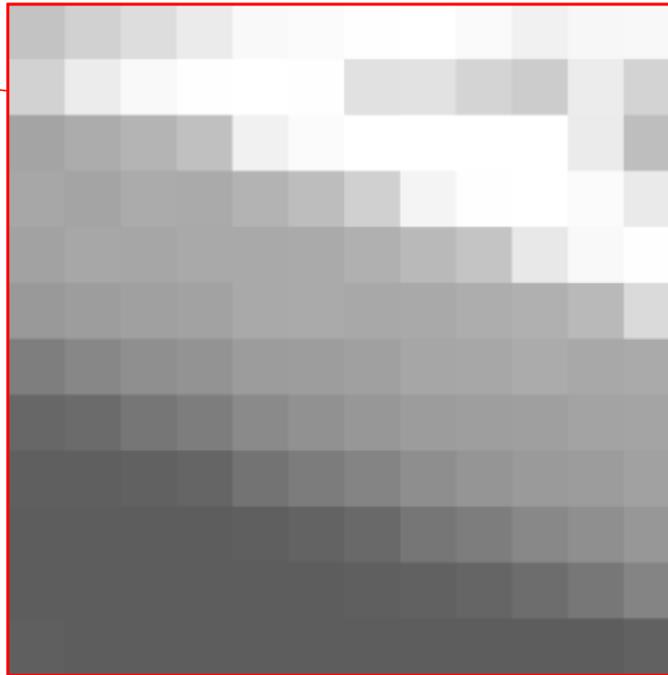
(از روی مثلث‌های متشابه)
 مقیاس / فاصله نامعین است.

تصویر



تصویر

مقادیر پیکسل‌ها



195	209	221	235	249	251	254	255	250	241	247	248
210	236	249	254	255	254	225	226	212	204	236	211
164	172	180	192	241	251	255	255	255	255	235	190
167	164	171	170	179	189	208	244	254	255	251	234
162	167	166	169	169	170	176	185	196	232	249	254
153	157	160	162	169	170	168	169	171	176	185	218
126	135	143	147	156	157	160	166	167	171	168	170
103	107	118	125	133	145	151	156	158	159	163	164
095	095	097	101	115	124	132	142	117	122	124	161
093	093	093	093	095	099	105	118	125	135	143	119
093	093	093	093	093	093	095	097	101	109	119	132
095	093	093	093	093	093	093	093	093	093	093	119

pixel = picture + element

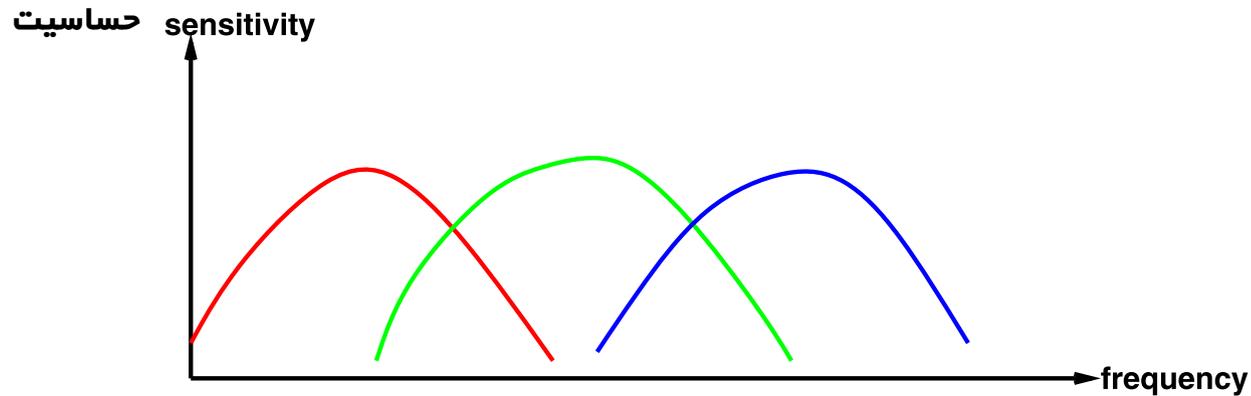
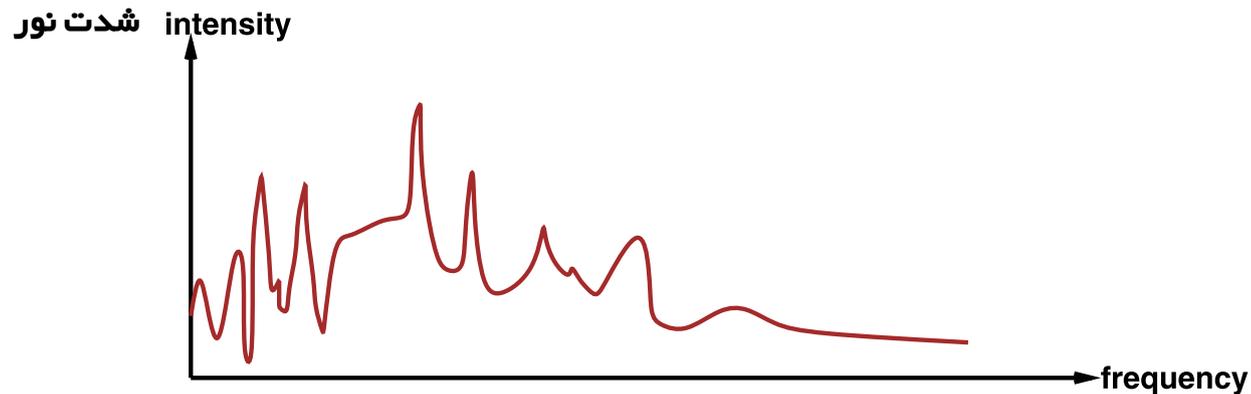
$I(x, y, t)$ is the intensity at (x, y) at time t

دوربین CCD ≈ 12 میلیون پیکسل
چشم انسان ≈ 240 میلیون پیکسل، 0.25 ترابیت بر ثانیه

بینایی رنگی

COLOR VISION

شدت نور با فرکانس تغییر می کند \Leftarrow سیگنال بی نهایت بعدی



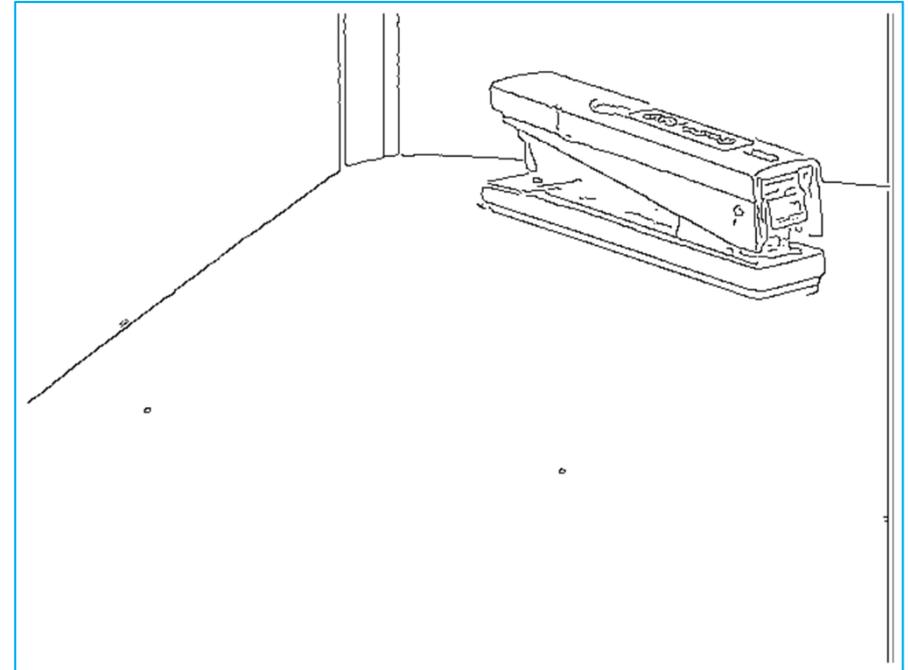
چشم انسان سه نوع سلول حساس به رنگ دارد؛
هر یک از کل سیگنال بی نهایت بعدی انتگرال می گیرد \Leftarrow بردار شدت نور سه عنصری

ادراک: بینایی ماشینی

۲

عملیات
پردازش
تصویر
اولیه

آشکار سازی لبه

EDGE DETECTION

Edges in image \Leftarrow discontinuities in scene:

- 1) depth
- 2) surface orientation
- 3) reflectance (surface markings)
- 4) illumination (shadows, etc.)

لبه‌ها در تصویر \Rightarrow ناپیوستگی‌ها در صحنه:

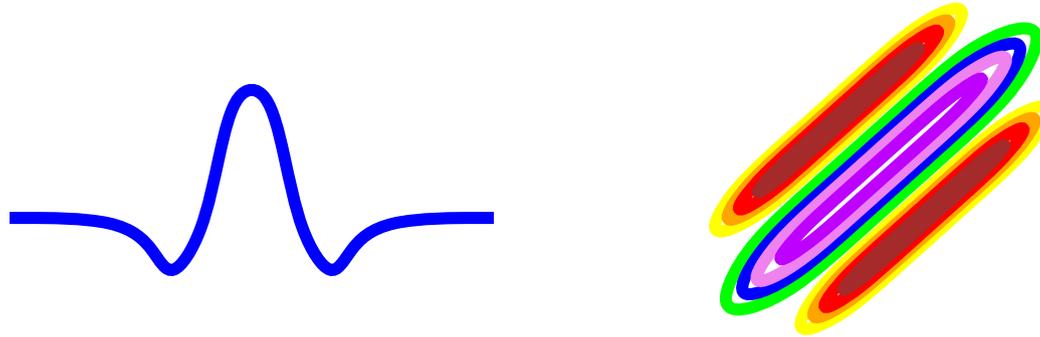
- (۱) عمق
- (۲) جهت رویه
- (۳) بازتاب (علامت‌گذاری‌های رویه)
- (۴) تابش (سایه‌ها و ...)

آشکار سازی لبه

الگوریتم

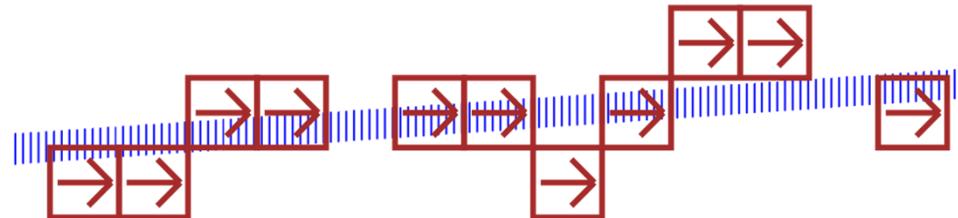
EDGE DETECTION

۱) تصویر را با فیلترهای مایل در حوزه‌ی مکان (احتمالاً multiscale) کانوالو کنید.



۲) پیکسل‌های بالای مقدار آستانه (above-threshold) را با جهت لبه برچسب‌گذاری کنید.

۳) پاره‌خط‌های «تمیز» را با ترکیب پیکسل‌های لبه‌ی متوالی با جهت یکسان استنتاج کنید.



ادراک: بینایی ماشینی

۳

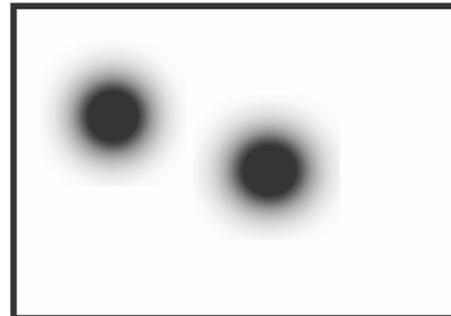
بازشناسی
اشیا
از روی
ظاهر

بازشناسی اشیا از روی ظاهر

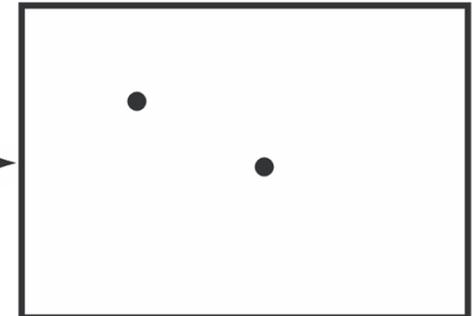
مثال



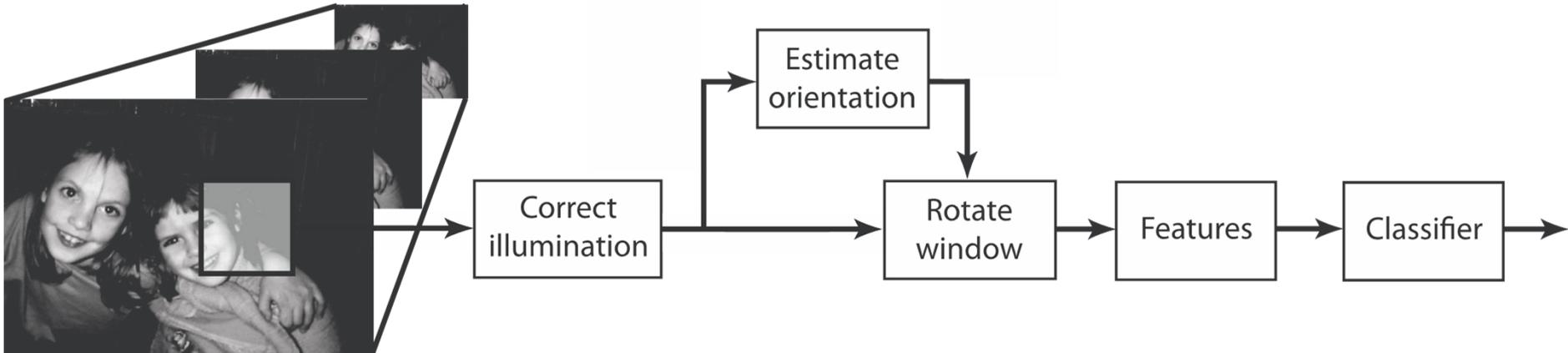
Image



Responses

Non-maximal
suppression

Detections



ادراک: بینایی ماشینی

۴

بازسازی
دنیای
سه بعدی

بازسازی دنیای سه بعدی

اشاره‌هایی از دانایی پیشین

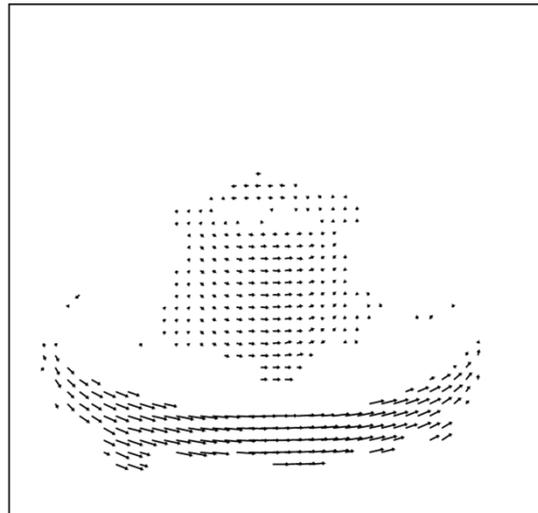
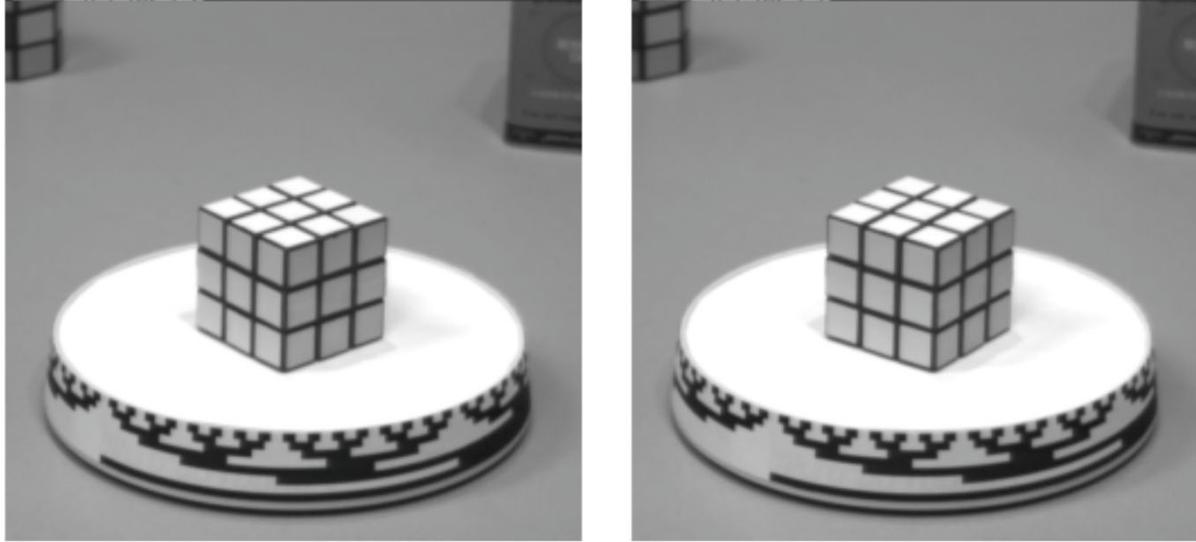
RECONSTRUCTING THE 3D WORLD: CUES FROM PRIOR KNOWLEDGE

فرض می‌کند ... <i>Assumes ...</i>	شکل به دست آمده از ... <i>Shape From ...</i>
بدنه‌های صلب، حرکت پیوسته <i>rigid bodies, continuous motion</i>	حرکت <i>Motion</i>
جامد، همجوار، بدنه‌های غیرتکراری <i>solid, contiguous, non-repeating bodies</i>	استریو <i>Stereo</i>
بافت یکنواخت <i>uniform texture</i>	بافت <i>Texture</i>
بازتاب یکنواخت <i>uniform reflectance</i>	سایه زنی <i>Shading</i>
انحنای می‌نیم <i>minimum curvature</i>	کانتور <i>Contour</i>

بازسازی دنیای سه بعدی

حرکت

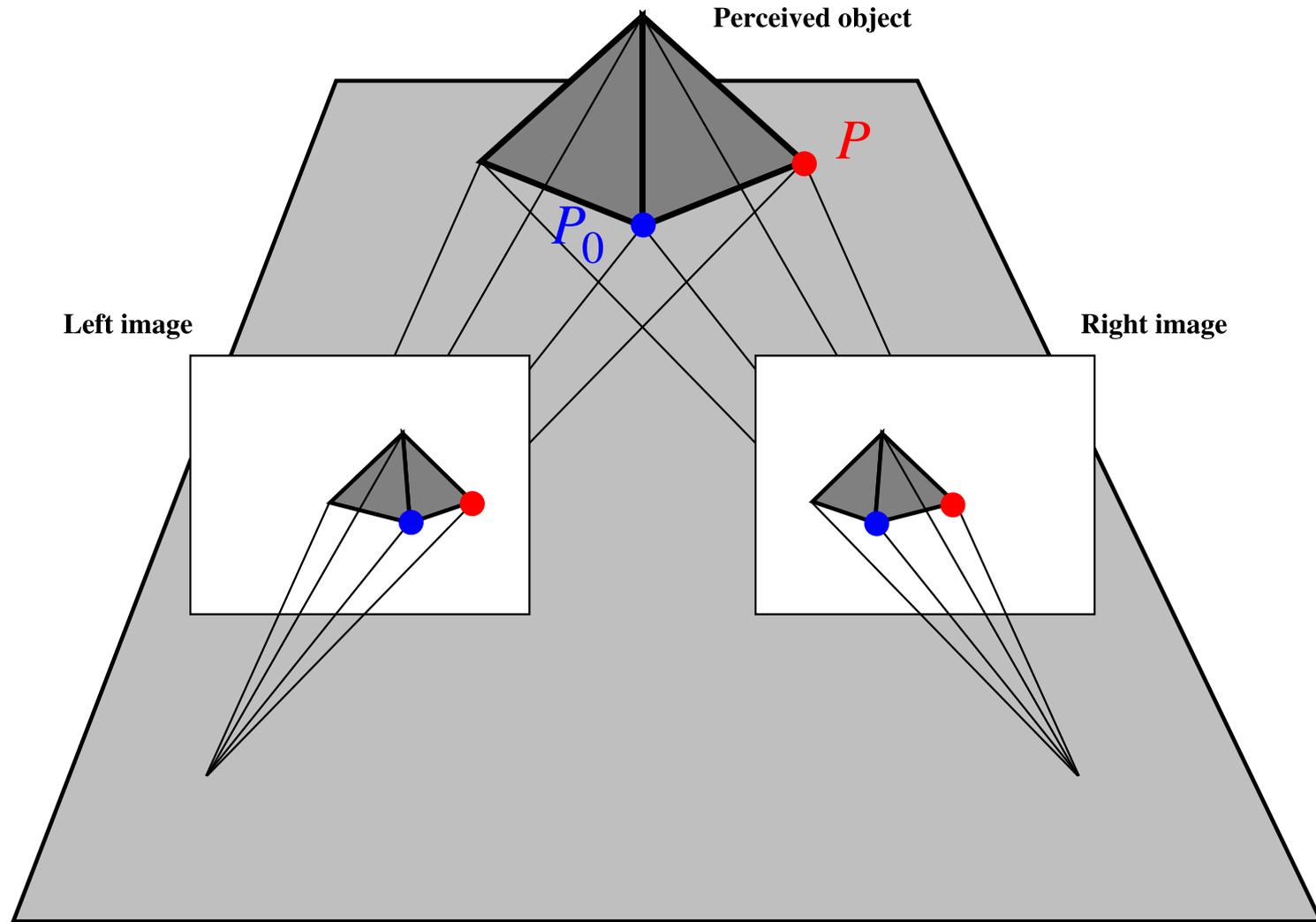
MOTION



بازسازی دنیای سه بعدی

استریو

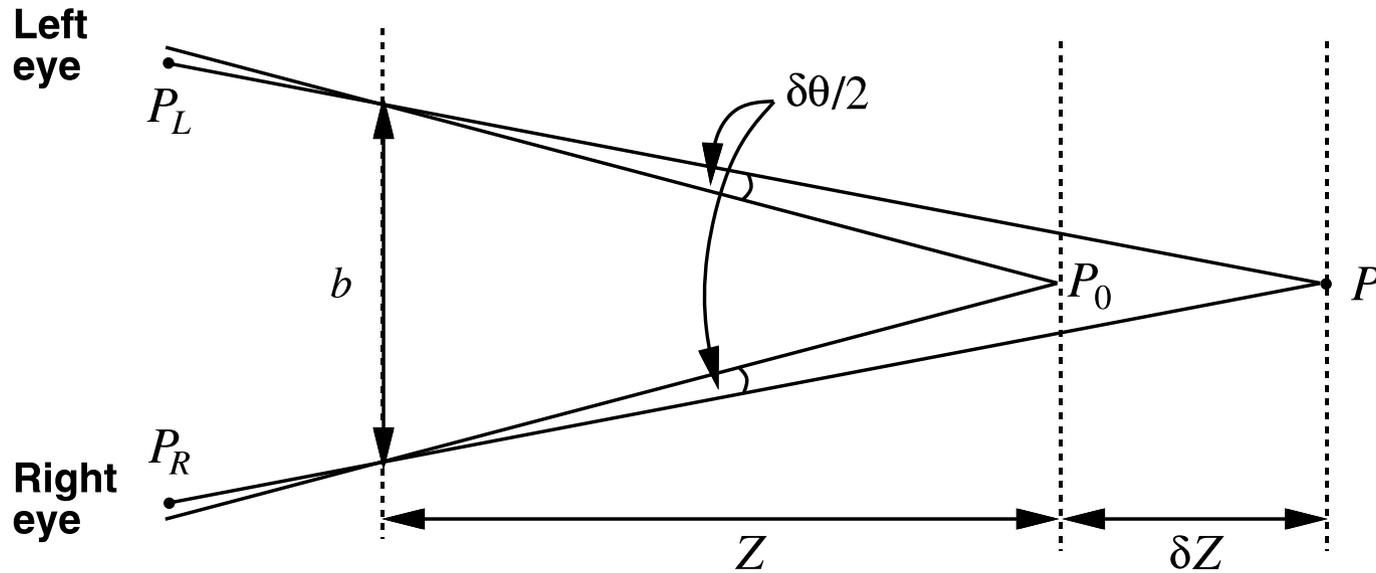
STEREO



بازسازی دنیای سه‌بعدی

استریو: تعیین عمق

STEREO



Simple geometry: $\delta Z = Z^2 \delta \theta / (-b)$

Physiology: $\delta \theta \geq 2.42 \times 10^{-5}$ radians, $b = 6\text{cm}$

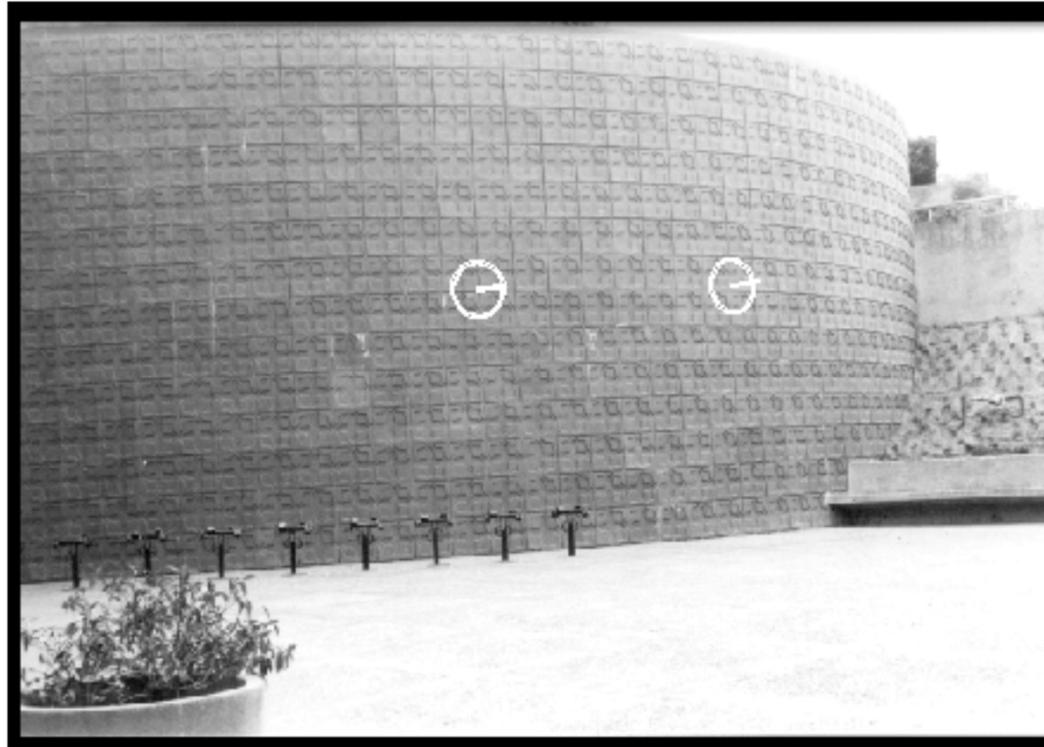
$Z = 30\text{cm} \Rightarrow \delta Z \approx 0.04\text{mm}$

$Z = 30\text{m} \Rightarrow \delta Z \approx 40\text{cm}$

Large baseline \Rightarrow better resolution!

بازسازی دنیای سه بعدی

بافت

TEXTURE

ایده: فرض می شود بافت واقعی یکنواخت است،
شکل رویه ای که این اعوجاج را ایجاد کرده است، محاسبه کنید.

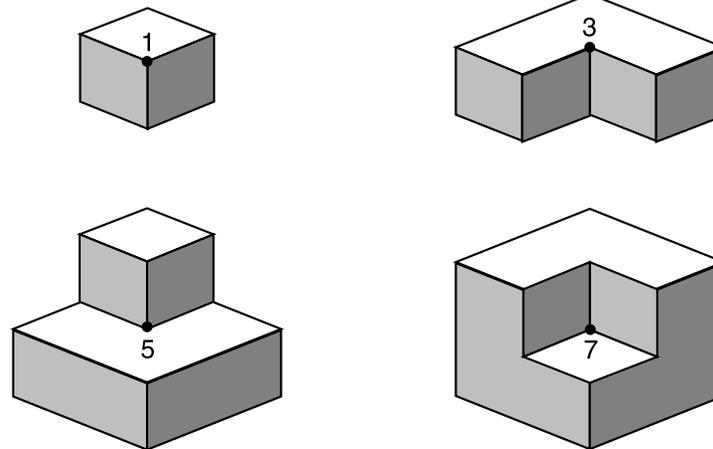
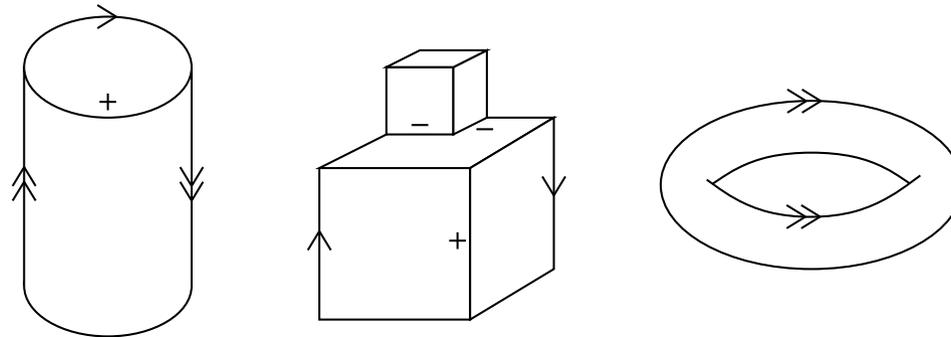
ایده ی مشابه برای سایه هم کار می کند:

فرض می کنیم بازتاب و ... یکنواخت باشد، اما عدم بازتاب ها محاسبات غیرمحلی شدت نور دریافت شده را موجب می شوند.
← گودی ها کم عمق تر از آنچه هستند به نظر می رسد.

بازسازی دنیای سه بعدی

انواع لبه‌ها و رأس‌ها

EDGE AND VERTEX TYPES

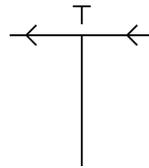
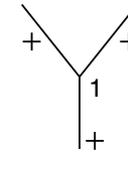
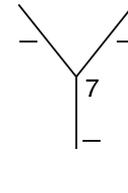
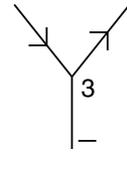
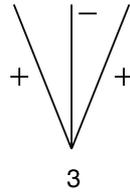
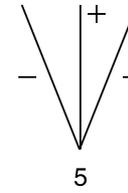
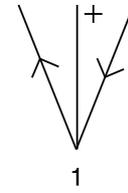
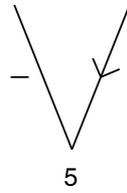
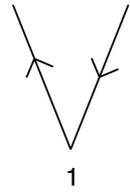
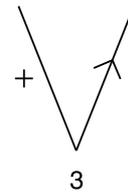
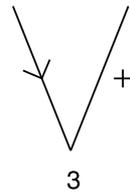
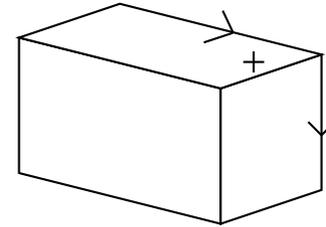
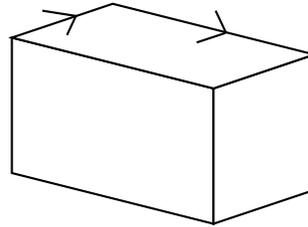
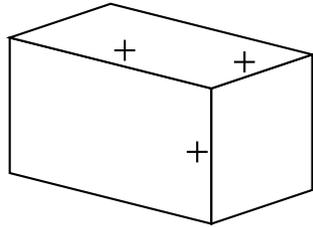


دنیای اشیای چندوجهی با رئوس سه وجهی را در نظر می‌گیریم.

بازسازی دنیای سه بعدی

برچسب‌های رأس / لبه

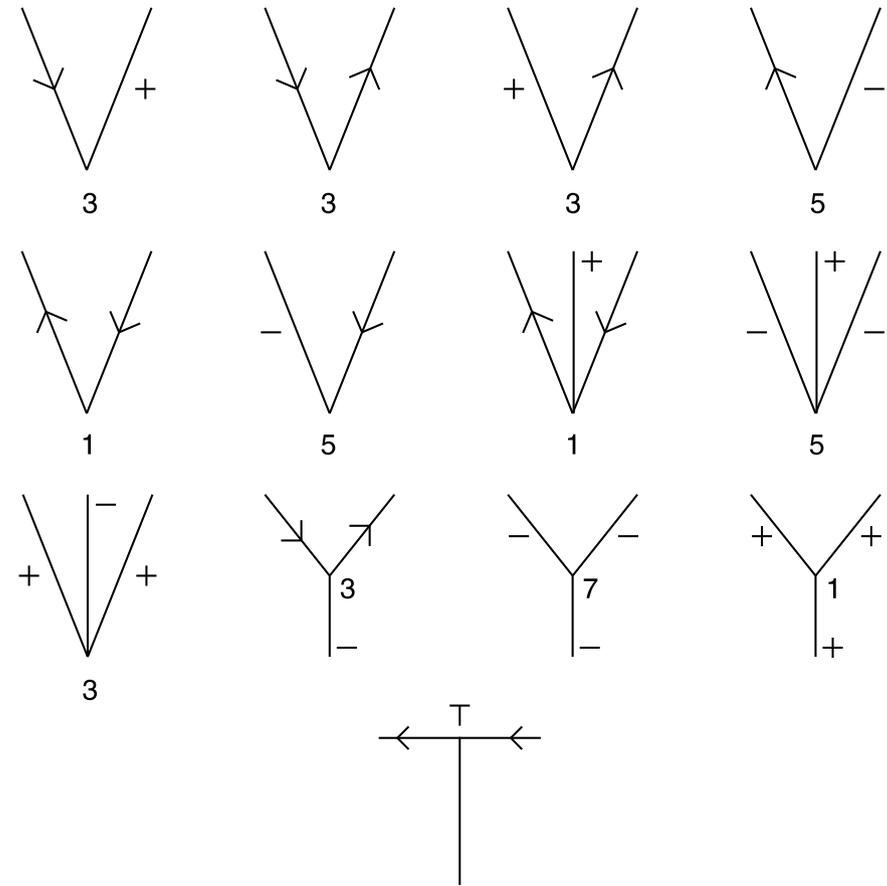
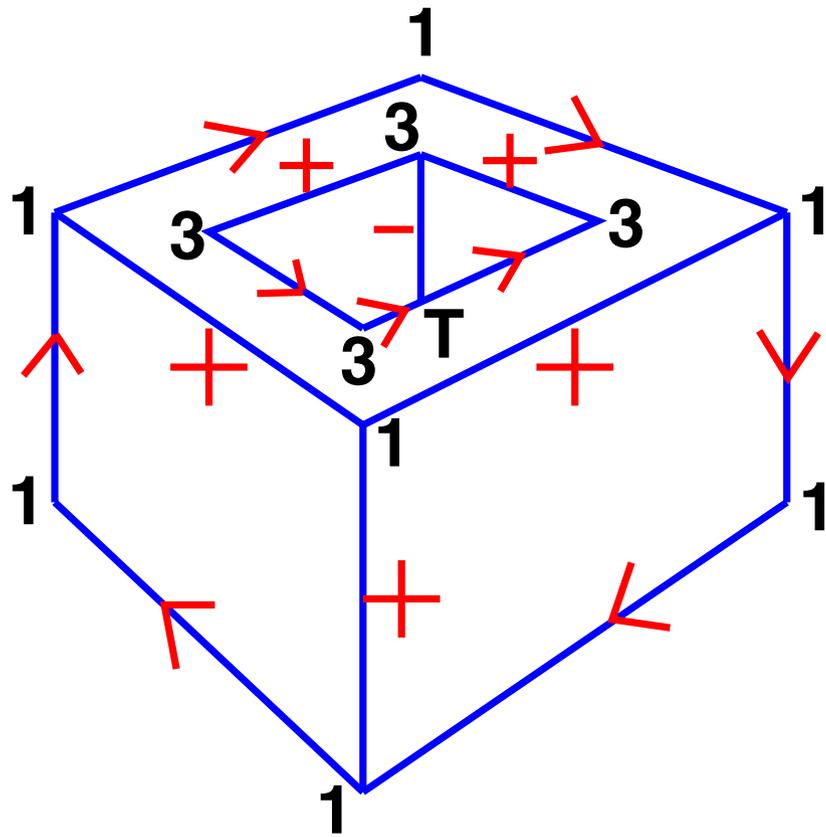
VERTEX/EDGE LABELS



بازسازی دنیای سه بعدی

برچسب زنی رأس / لبه : مثال

VERTEX/EDGE LABELLING: EXAMPLE



مسئله‌ی برچسب زنی رأس / لبه به عنوان یک مسئله‌ی ارضای قید (CSP):

متغیرها = لبه‌ها

قیدها = پیکربندی‌های ممکن گره‌ها

۵

بازشناسی
اشیا
از روی
اطلاعات
ساختاری

بازشناسی اشیا

OBJECT RECOGNITION

ایده‌ی ساده برای بازشناسی اشیا:
 - استخراج شکل‌های سه‌بعدی از تصویر
 - تطابق در برابر «کتابخانه‌ی شکل‌ها»

- استخراج رویه‌های انحنادار از تصویر
- بازنمایی شکل شیئی استخراج شده
- بازنمایی شکل و تغییرپذیری طبقه‌های اشیا در کتابخانه
- قطعه‌بندی نامناسب، پنهان شدن (قرار گرفتن شیئی روی دیگری)
- نامعلوم بودن نورپردازی، سایه‌ها، علامت‌گذاری‌ها، نویز، پیچیدگی و ...

مسائل بازشناسی شیئی

- نمایه‌سازی کتابخانه با اندازه‌گیری خصوصیات نامتغیر اشیا
- ترازبندی ویژگی‌های تصویر با ویژگی‌های شیئی پروجکت شده در کتابخانه
- تطابق تصویر در برابر چندین نمای ذخیره شده (جنبه‌ها) از کتابخانه‌ی اشیا
- روش‌های یادگیری ماشینی بر اساس آماره‌های تصویر

روی‌کردها

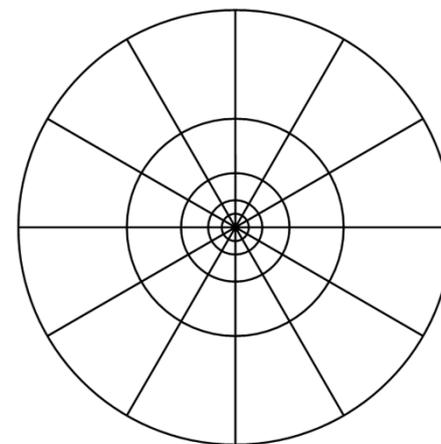
بازشناسی اشیا

روش تطابق زمینه‌ی شکل (۱ از ۳)

SHAPE-CONTEXT MATCHING

ایده‌ی پایه:

شکل (یک مفهوم رابطه‌ای) را به مجموعه‌ای ثابت از **خصیصه‌ها** تبدیل کنید؛
با استفاده از **بستر مکانی** هر یک از مجموعه نقاط ثابت روی رویه‌ی شکل.

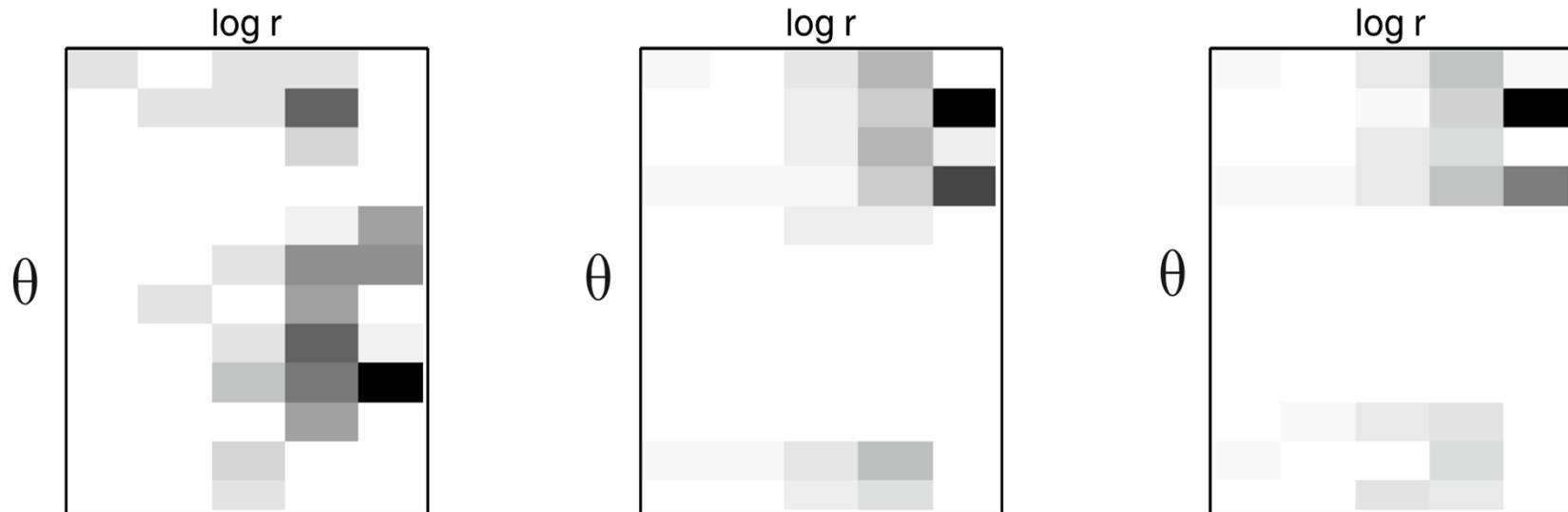


بازشناسی اشیا

روش تطابق زمینه‌ی شکل (۲ از ۳)

SHAPE-CONTEXT MATCHING

هر نقطه با هیستوگرام بستر محلی خودش توصیف می‌شود.
(تعداد نقاطی که در هر بسته‌ی مشبکه‌ی نمودار لگاریتمی-قطبی می‌افتد)



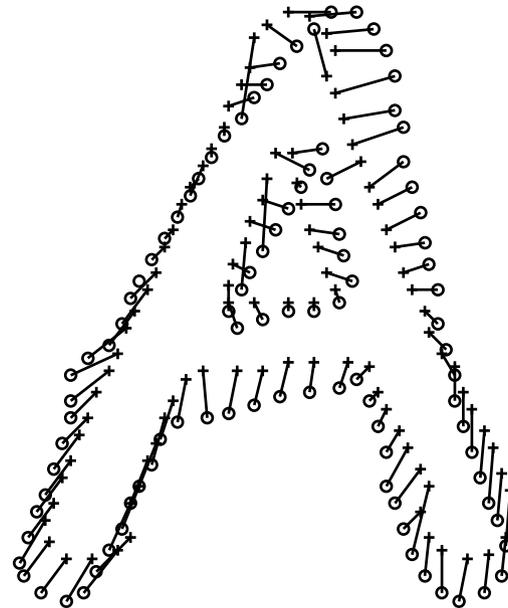
بازشناسی اشیا

روش تطابق زمینه‌ی شکل (۳ از ۳)

SHAPE-CONTEXT MATCHING

فاصله‌ی کل بین شکل‌ها

با مجموع فاصله‌های میان نقاط متناظر تحت بهترین تطابق تعیین می‌شود.



الگوریتم یادگیری ساده‌ی نزدیک‌ترین همسایه، 0.63% خطا بر روی مجموعه داده اعداد NIST نشان می‌دهد.

ادراک: بینایی ماشینی

۶

استفاده از
بینایی

استفاده از بینایی کامپیوتری

کاربردها

بازشناسی
پلاک خودرو

*Car Plate
Recognition*

بازشناسی
چهره

*Face
Recognition*

آشکارسازی
چهره

*Face
Detection*

بازیابی
تصویر مبتنی
بر محتوا

*Content-Based
Image Retrieval*

شناسایی
زیست‌سنجی

*Biometric
Identification*

بازشناسی
نویسه‌ها

*Character
Recognition
(OCR)*

بازشناسی
دست‌خط

*Handwriting
Recognition*

بازشناسی
امضا

*Signature
Recognition*

واسط کاربر
مبتنی بر
بینایی

*Vision-Based
User Interface*

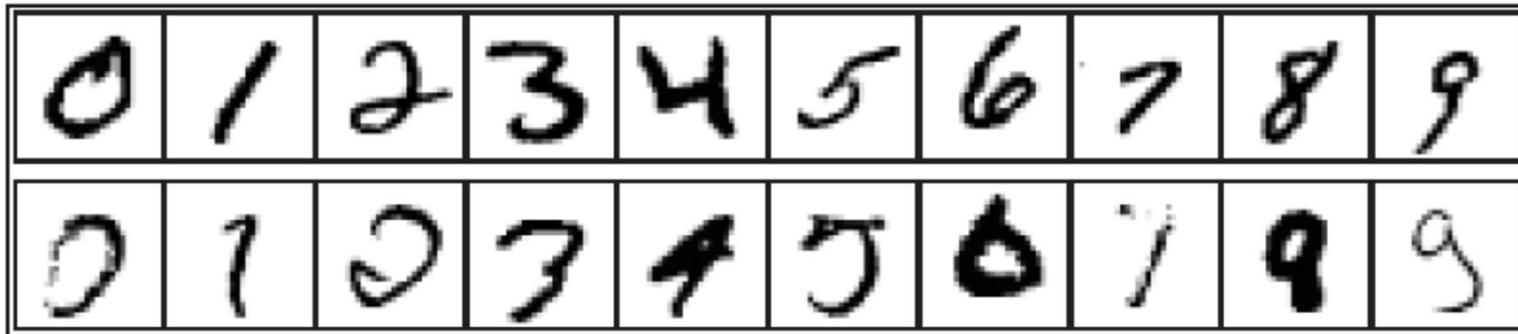
حفاظت مبتنی
بر تصویر

*Image-Based
Protection*

بازشناسی اعداد دستنویس

مثالی از بازشناسی اشیا

HANDWRITTEN DIGIT RECOGNITION



3-nearest-neighbor = 2.4% error

400–300–10 unit MLP = 1.6% error

LeNet: 768–192–30–10 unit MLP = 0.9% error

Current best (kernel machines, vision algorithms) \approx 0.6% error

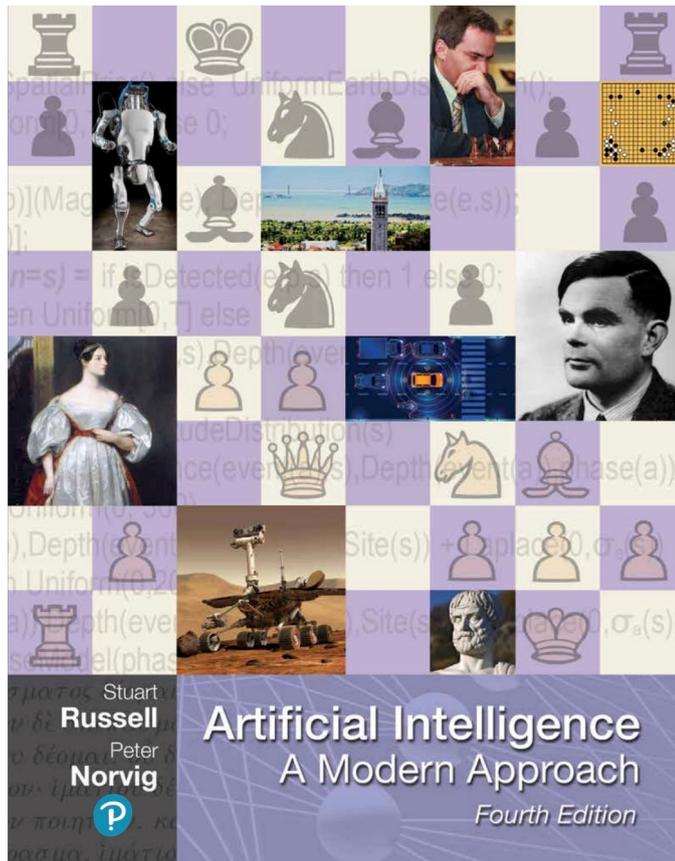
هوش مصنوعی

ادراک: بینایی ماشینی

۷

منابع،
مطالعه،
تکلیف

منبع اصلی



Stuart Russell and Peter Norvig,
Artificial Intelligence: A Modern Approach,
 4th Edition, Prentice Hall, 2020.

Chapter 25

CHAPTER 25

COMPUTER VISION

In which we connect the computer to the raw, unwashed world through the eyes of a camera.

Most animals have eyes, often at significant cost: eyes take up a lot of space; use energy; and are quite fragile. This cost is justified by the immense value that eyes provide. An agent that can see can predict the future—it can tell what it might bump into; it can tell whether to attack or to flee or to court; it can guess whether the ground ahead is swampy or firm; and it can tell how far away the fruit is. In this chapter, we describe how to recover information from the flood of data that comes from eyes or cameras.

25.1 Introduction

Vision is a perceptual channel that accepts a **stimulus** and reports some representation of the world. Most agents that use vision use **passive sensing**—they do not need to send out light to see. In contrast, **active sensing** involves sending out a signal such as radar or ultrasound, and sensing a reflection. Examples of agents that use active sensing include bats (ultrasound), dolphins (sound), abyssal fishes (light), and some robots (light, sound, radar). To understand a perceptual channel, one must study both the physical and statistical phenomena that occur in sensing and what the perceptual process should produce. We concentrate on vision in this chapter, but robots in the real world use a variety of sensors to perceive sound, touch, distance, temperature, global position, and acceleration.

A **feature** is a number obtained by applying simple computations to an image. Very useful information can be obtained directly from features. The wumpus agent had five sensors, each of which extracted a single bit of information. These bits, which are features, could be interpreted directly by the program. As another example, many flying animals compute a simple feature that gives a good estimate of time to contact with a nearby object; this feature can be passed directly to muscles that control steering or wings, allowing very fast changes of direction. This **feature extraction** approach emphasizes simple, direct computations applied to sensor responses.

The **model-based** approach to vision uses two kinds of models. An **object model** could be the kind of precise geometric model produced by computer aided design systems. It could also be a vague statement about general properties of objects, for example, the claim that all faces viewed in low resolution look approximately the same. A **rendering model** describes the physical, geometric, and statistical processes that produce the stimulus from the world. While rendering models are now sophisticated and exact, the stimulus is usually ambiguous. A white object under low light may look like a black object under intense light. A small, nearby object may look the same as a large, distant object. Without additional evidence,

Feature