

Intro to Health Data Science

MMCi Block 1

Matthew Engelhard

Introductions!

your name

+ a brief note on your background and role

+ which of the 5 topics interests you most

(intro to DS, model learning, computer vision, natural language processing, sequential data)

Course Overview

We will learn about state-of-the-art data science techniques that are now beginning to impact clinical practice.

- How are these techniques different from what has come before?
- How are they the same?
- And what do you need to know to take advantage of this tech?

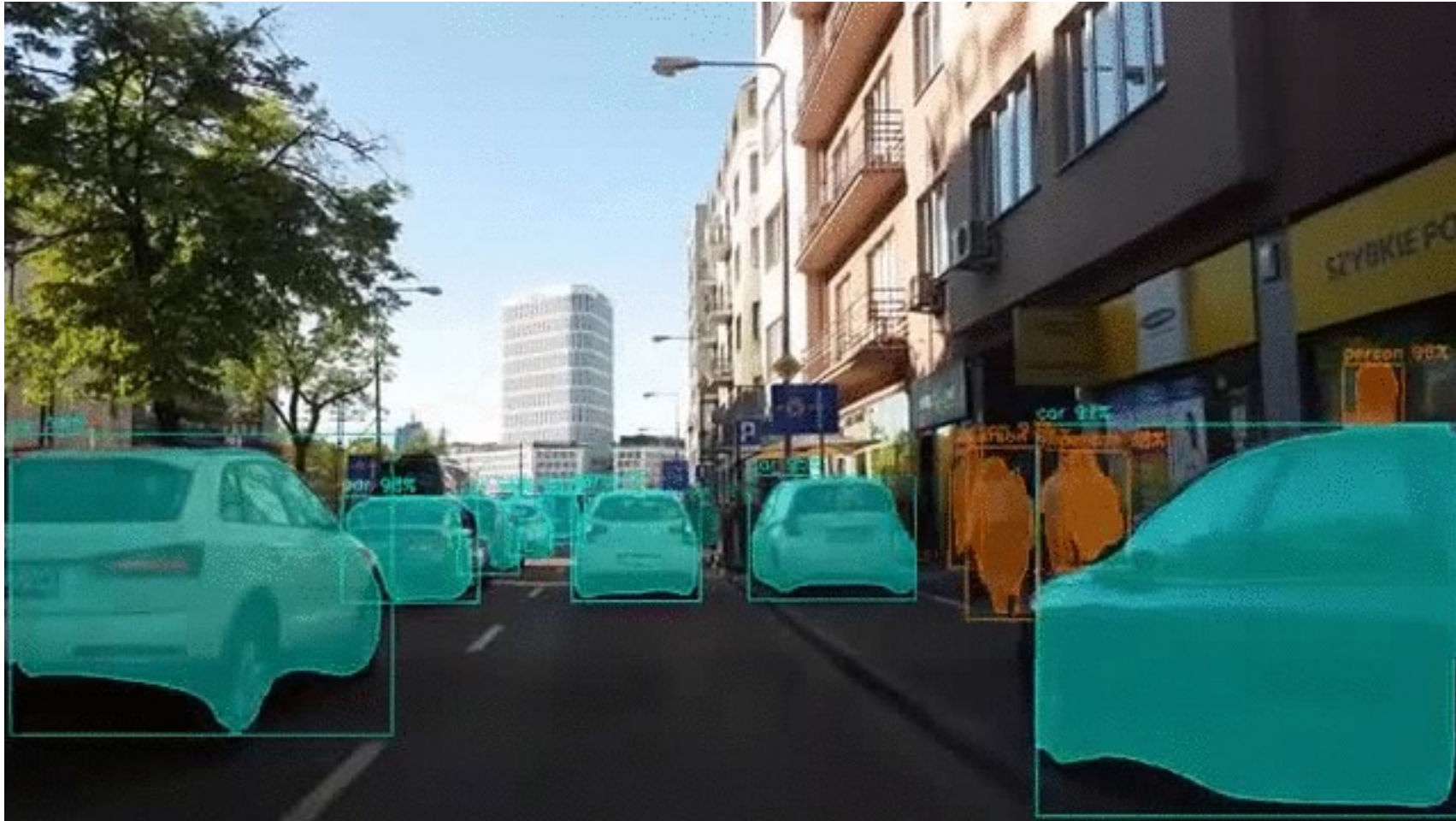
I know that most of you are NOT going to be data scientists.

But you *will* work with data scientists, and you *will* have to make decisions about what models to use and how to use them. It is important to know enough to get in the weeds with the data scientists, because if applied/evaluated incorrectly, these models are certain to be unhelpful and *likely to be harmful*.

A Brief Tour of DS in 2021

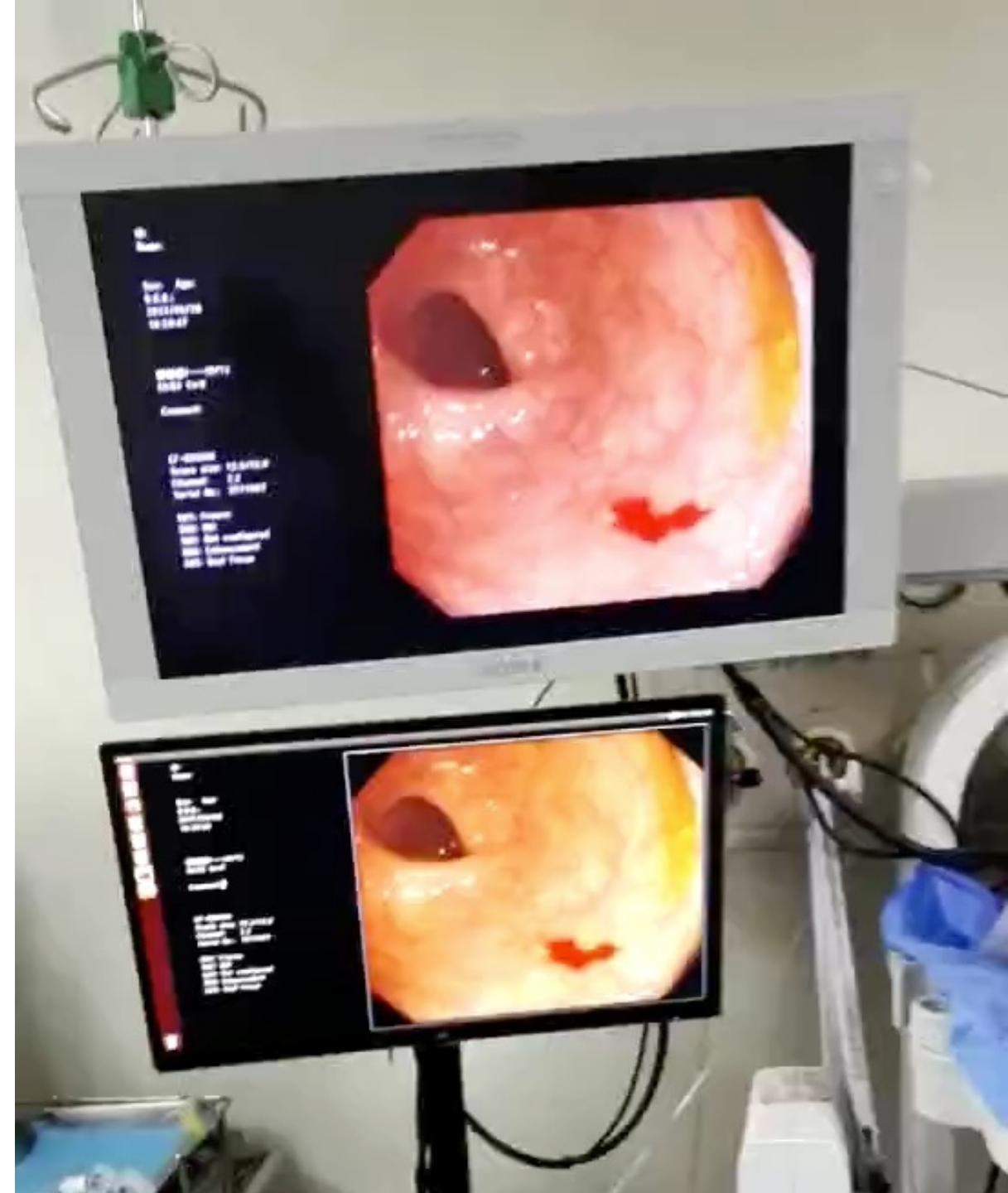
General Application -> Biomedical Application

Computer Vision: Real-Time Object Detection (& segmentation)



Real-time polyp
detection:

a second set of eyes
during colonoscopy



Natural Language Processing: Text Generation

SYSTEM PROMPT (HUMAN-WRITTEN)

A train carriage containing controlled nuclear materials was stolen in Cincinnati today. Its whereabouts are unknown.

“Better Language Models and Their Implications”

2/14/19

OPENAI

MODEL COMPLETION (MACHINE-WRITTEN, FIRST TRY)

The incident occurred on the downtown train line, which runs from Covington and Ashland stations.

In an email to Ohio news outlets, the U.S. Department of Energy said it is working with the Federal Railroad Administration to find the thief.

“The theft of this nuclear material will have significant negative consequences on public and environmental health, our workforce and the economy of our nation,” said Tom Hicks, the U.S. Energy Secretary, in a statement. “Our top priority is to secure the theft and ensure it doesn’t happen again.”

The stolen material was taken from the University of Cincinnati’s Research Triangle Park nuclear research site, according to a news release from Department officials. The Nuclear Regulatory Commission did not immediately release any information.

According to the release, the U.S. Department of Energy’s Office of Nuclear Material Safety and Security is leading that team’s investigation.

“The safety of people, the environment and the nation’s nuclear stockpile is our highest priority,” Hicks said. “We will get to the bottom of this and make no excuses.

Natural Language Processing: Find Predictive Text in Notes

Passage (from note) | Change in predicted autism dx log-odds

subjective intake chief complaint problems with sleep, inattention, and behavioral concerns both in the home and school setting. DATE, recently more anger and recent tic like behavior +6.95

psychologist presenting problem NAME is a 3 year, 4 month old female who was referred for a neurodevelopmental assessment due to concerns regarding her overall development, behavior, and social emotional functioning and to assess for autism spectrum disorder +6.82

problem list diagnosis • disruptive behavior disorder • impaired speech articulation • daytime enuresis • other subjective visual disturbances • hypermetropia of both eyes • adhd attention deficit +6.81

problem list diagnosis • anemia of prematurity • history of colitis • meconium tox for thc • extreme immaturity of newborn, 27 completed weeks • nasal congestion of newborn • presumed +6.78

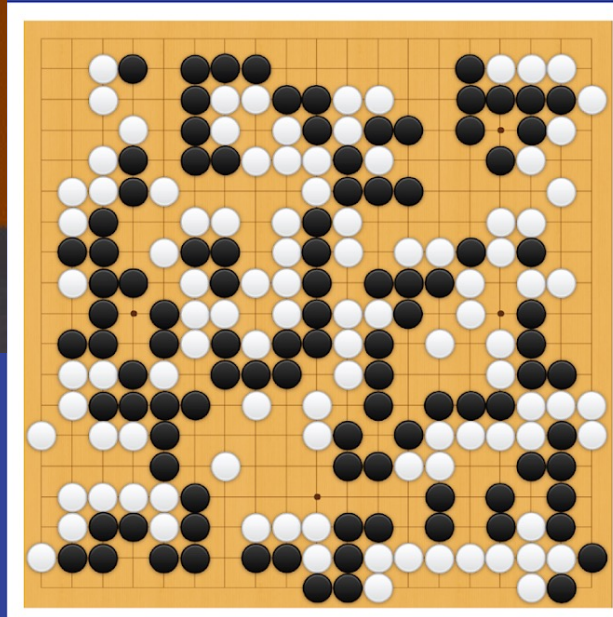
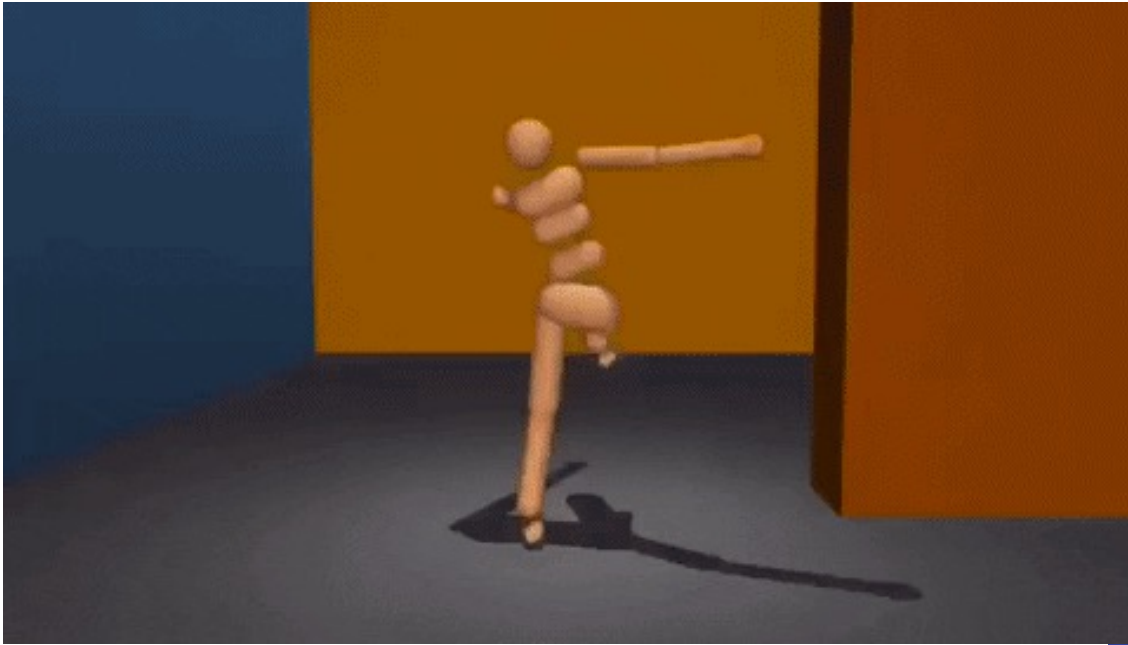
motor delay DATE • hypotonia DATE • clasped thumb DATE • polydactyly DATE • developmental +6.74

therapy NAME was seen for developmental support during rop eye exam today. the +6.65

← Developmental and behavioral concerns are highly predictive

← Premature birth and perinatal complications are also highly predictive

Reinforcement Learning: goal-directed sequential decision-making



THE ULTIMATE GO CHALLENGE

GAME 3 OF 3

27 MAY 2017



vs



AlphaGo

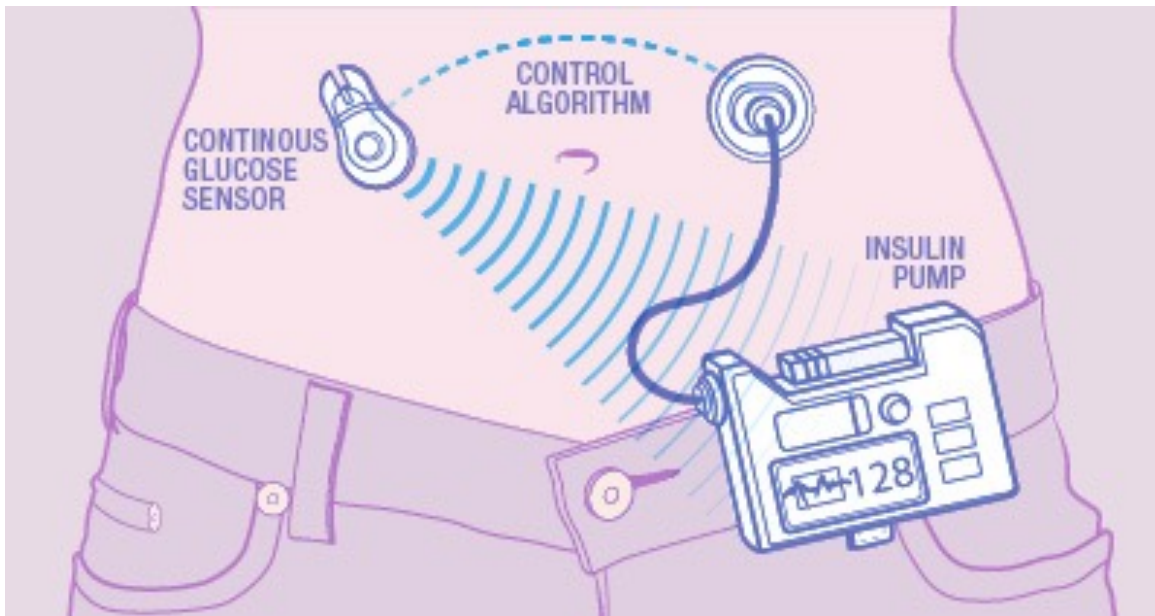
Winner of Match 3

Ke Jie

RESULT B + Res

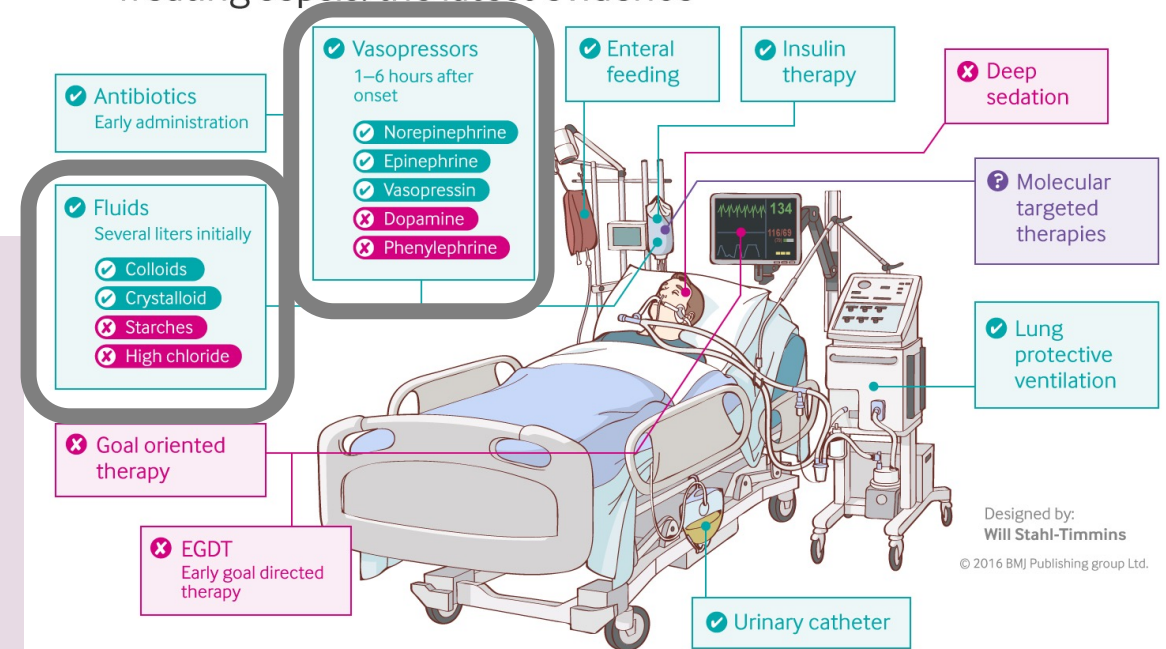
Reinforcement Learning in Medicine

Closed-loop blood glucose control ("artificial pancreas")



<https://www.mayo.edu/research/labs/artificial-pancreas/overview>

Treating sepsis: the latest evidence



Fluid and vasopressor administration for sepsis treatment

Gotts JE, Matthay MA. Sepsis: pathophysiology and clinical management. *bmj*. 2016 May 23;353(i1585).

The Current DS Moment

Looking back to 2012...

Deep learning leapt forward in '12 and beat humans in '15

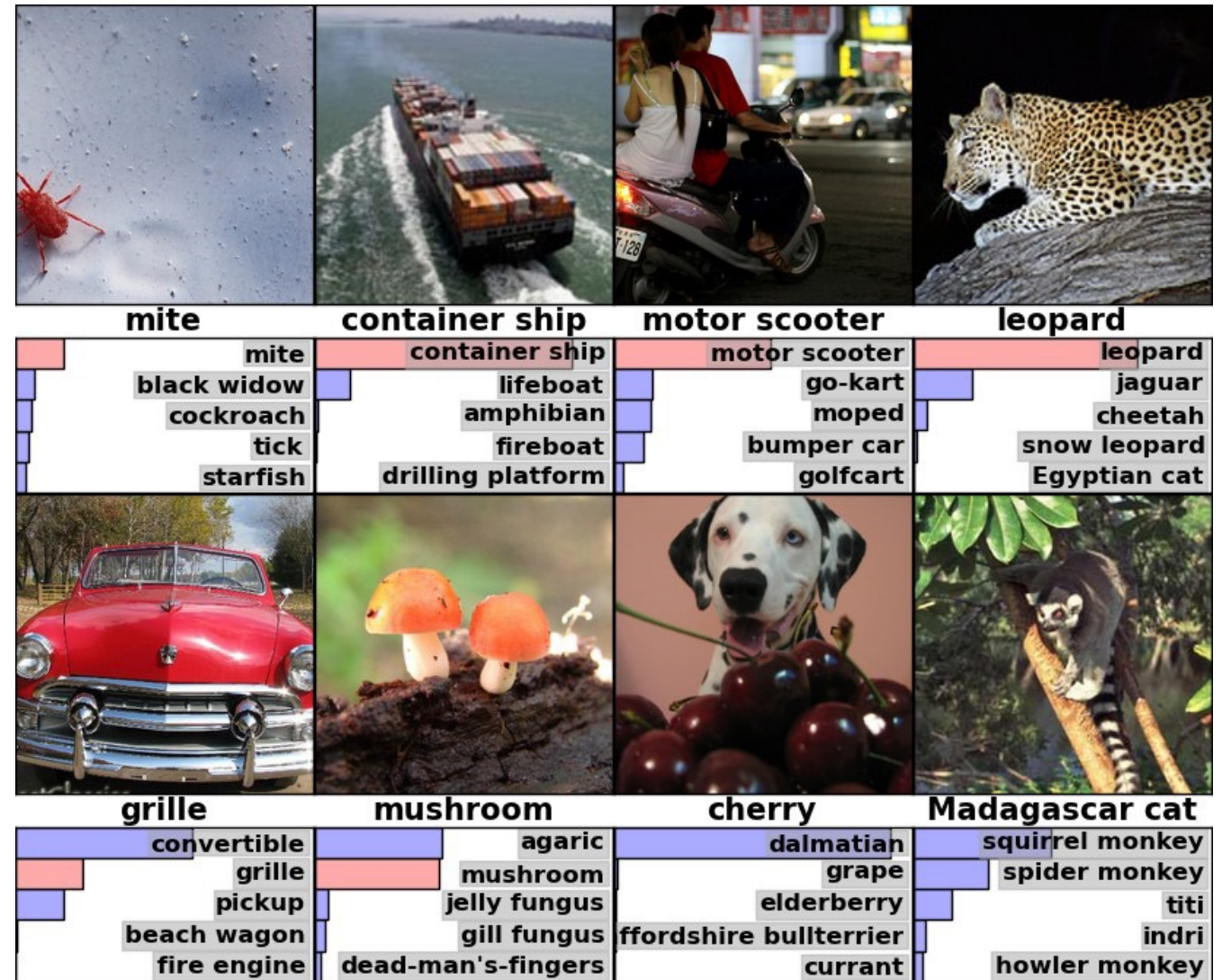
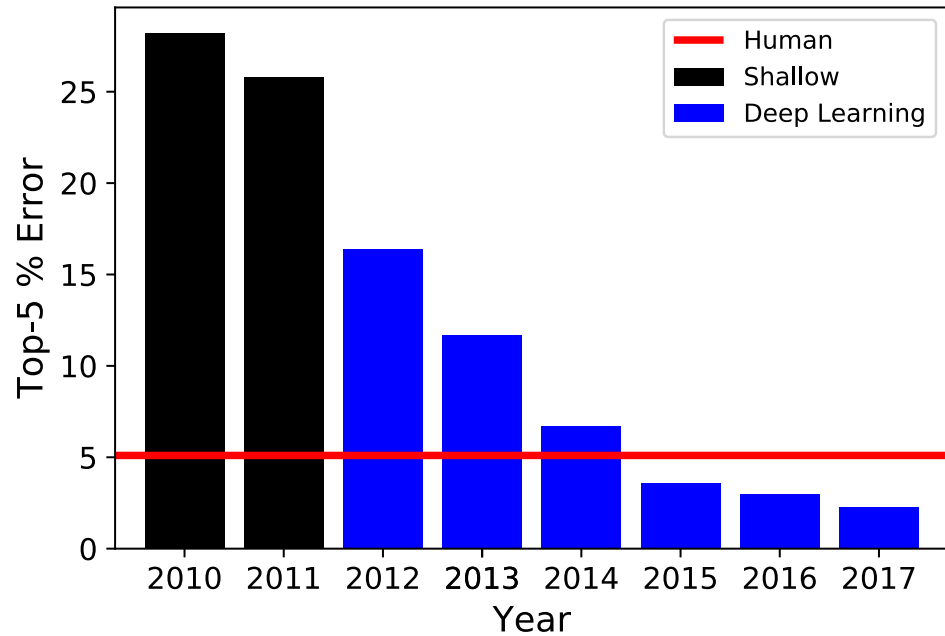


Figure from Krizhevsky et al 2012

Deep Learning now surpasses humans in a variety of tasks



Most recently, DL has surpassed humans in language tasks

Microorganisms or toxins that successfully enter an organism encounter the cells and mechanisms of the innate immune system. The innate response is usually triggered when microbes are identified by pattern recognition receptors, which recognize components that are conserved in microorganisms, or when damaged, injured cells release signals, many of which (but not all) are recognized by those that recognize pathogens. Innate immunity means these systems respond to pathogens but do not confer long-lasting immunity against them. The innate immune system is the dominant system of host defense.

What part of the innate immune system identifies microbes and triggers immune response?

Ground Truth Answers: pattern recognition receptors receptors cells

Leaderboard

SQuAD2.0 tests the ability of a system to not only answer reading comprehension questions, but also abstain when presented with a question that cannot be answered based on the provided paragraph. How will your system compare to humans on this task?

Rank	Model	EM	F1
	Human Performance Stanford University (Rajpurkar & Jia et al. '18)	86.831	89.452
1 Mar 05, 2019	BERT + N-Gram Masking + Synthetic Self-Training (ensemble) Google AI Language https://github.com/google-research/bert	86.673	89.147
2 Mar 05, 2019	BERT + N-Gram Masking + Synthetic Self-Training (single model) Google AI Language https://github.com/google-research/bert	85.150	87.715

dominant system of defense?

innate immune

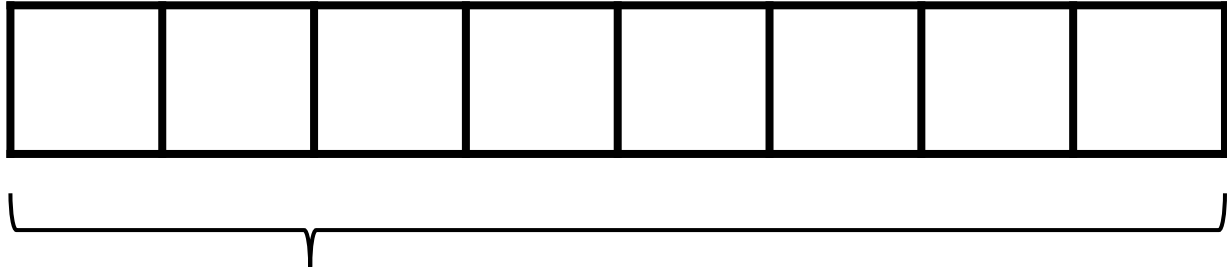
size components present in broad

microorganisms

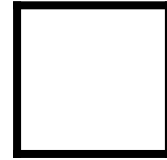
in a generic way, meaning it is

non-specific non-specific

All of these have, at their core, a predictive model



x , data/features for
a subject or patient



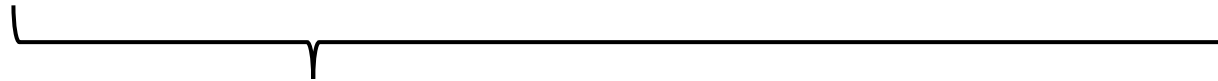
y , associated
value or label

End goal: predict y from x

Simple models often work well for clinical data!



Age
Pulse Rate
Mean BP
Temperature
Respiratory Rate
Hematocrit
WBC Count
Creatinine



x , data/features for
a subject or patient



Survival

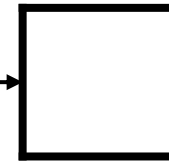
y , associated
value or label

End goal: predict odds
of hospital mortality

For complex data, more complex models are needed.



x , retinal image



y , referable diabetic
retinopathy

End goal: predict y from x

Course Objective

Understand of the capabilities and limitations of healthcare data science well enough to:

- (a) design and manage data science research and/or QA/QI projects
- (b) collaborate and communicate effectively with data scientists
- (c) critically evaluate data science models and methods, with an emphasis on rigorous model validation

Course Requirements and Materials

- [Let's take a look at the website](#)
- Questions & discussion about course requirements, materials, or activities
- Contact me: m.engelhard@duke.edu