

Day 9:

Speaker: Makarand Tapaswi, IIIT Hyderabad & Wadhawan AI

Title: Recent Advances in Video and Language Understanding.

Large improvements in computer vision

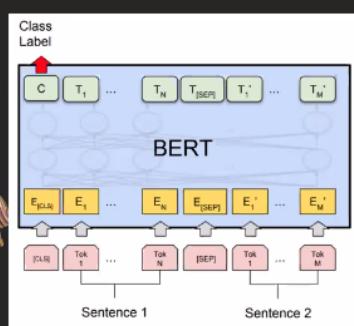


- Evolution in CV ↓
- Datasets with millions of annotated samples
 - Convolutional Neural Networks with 100+ million learnable parameters
 - Now, Transformers with billions of learnable parameters
 - Coupled with smart augmentation strategies for self-supervised training

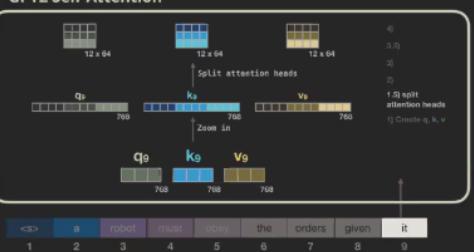
Intro – Representations – Stories – Instructional Videos

2

Large improvements in language representations



GPT2 Self-Attention



- Recent Dev's in Sequence Models
- Massive text corpora are harnessed to build strong language models
 - Examples: OpenAI's GPT and Google's BERT
 - Such models show transfer learning capabilities to numerous NLP tasks

Intro – Representations – Stories – Instructional Videos

3

Multimodal Learning

Joint understanding of Vision and Language

Image-text
Retrieval



Image
Captioning



A woman holding a banana up to her face.

Visual Question
Answering



Q. Is the woman smiling?
A. Yes (?)

How we perceive
↓

→ Yes

for some

by

No for
others

What does this picture show?



Objects and people:

- Mom
- Kid
- Book
- Teddy bear

What is happening?

- Mom is reading a book to her child

Raw
image

captioning

... but if it's a video?

The child
starts
crying after
closing the
book.



Rich information

- story in the book
- emotions of the parent and child
- why does the child cry?
- learn physics by observing that books on the shelf!

Example:

A multimodal audio-video-language stream



Intro – Representations – Stories – Instructional Videos

7

Real Life Problem

Joint understanding of Video and Language

Video Stories



Q. How did Jon Arryn die ?

A. Fever *

Representation learning



How to bake a cake

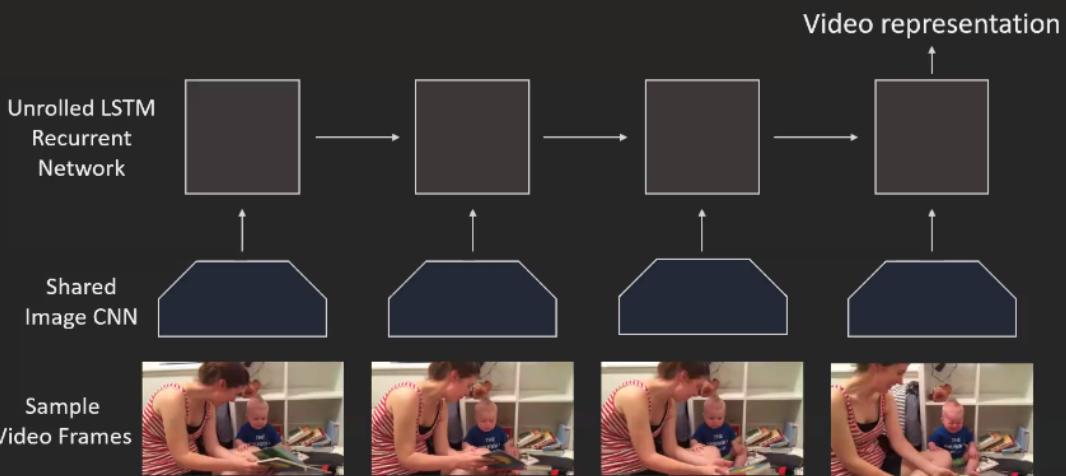
Intro – Representations – Stories – Instructional Videos

8

Whirlwind tour of
Video-Language Representations

9

Video: CNN-LSTM Encoders

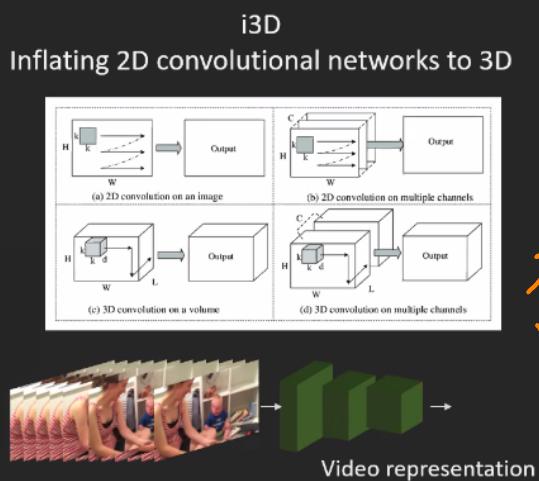


Reference: Long-term Recurrent Convolutional Networks for Visual Recognition and Description, PAMI 2014

Intro – Representations – Stories – Instructional Videos

10

Video: Time-aware Convolutional Networks



SlowFast Convolutional Networks



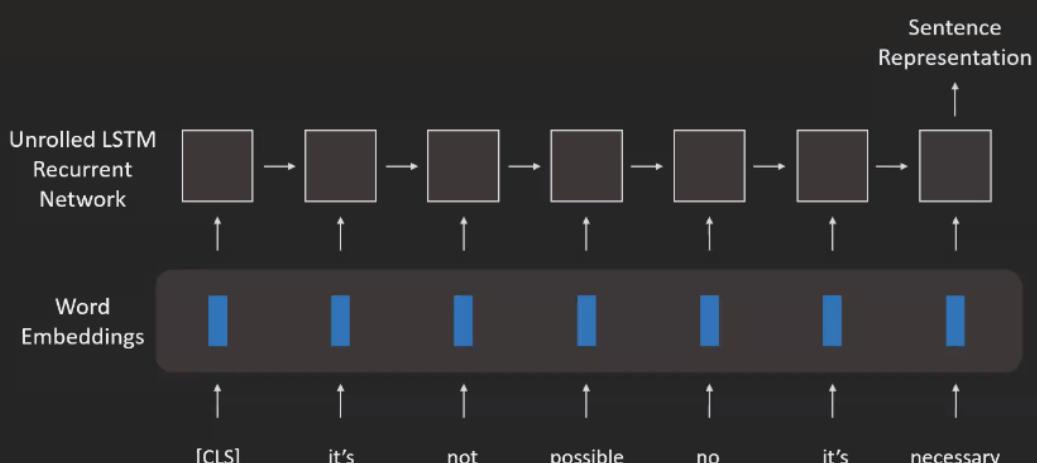
Reference: Quo Vadis, Action Recognition? A New Model and the Kinetics Dataset, CVPR 2017

Reference: SlowFast Networks for Video Recognition, ICCV 2019

Intro – Representations – Stories – Instructional Videos

11

Sentence: Word2Vec + RNN/LSTMs



Intro – Representations – Stories – Instructional Videos

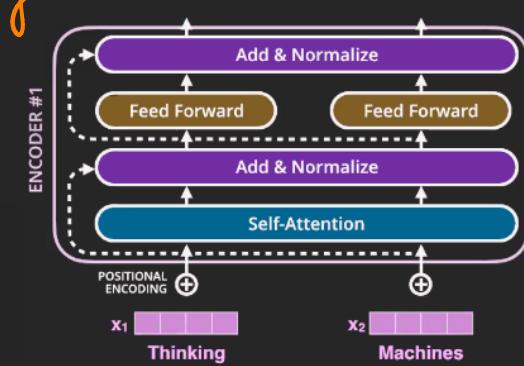
12

* Read → Attention is all you Need Paper! !

Transformers and the Self-Attention Layer

* key, query

* Self Attention



Input	Thinking	Machines
Embedding	x ₁ [purple]	x ₂ [purple]
Queries	q ₁ [green]	q ₂ [green]
Keys	k ₁ [blue]	k ₂ [blue]
Values	v ₁ [orange]	v ₂ [orange]
Score	q ₁ • k ₁ = 112	q ₁ • k ₂ = 96
Divide by 8 ($\sqrt{d_k}$)	14	12
Softmax	0.88	0.12
Softmax X Value	v ₁ [orange]	v ₂ [orange]
Sum	z ₁ [teal]	z ₂ [teal]

Reference: Attention is All you Need, NeurIPS 2017

Image Credits: The Illustrated Transformer, <https://jalammar.github.io/illustrated-transformer/>

Intro – Representations – Stories – Instructional Videos

13

BERT (very simplified diagram)!



Reference: BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding, arXiv 2018.10

Intro – Representations – Stories – Instructional Videos

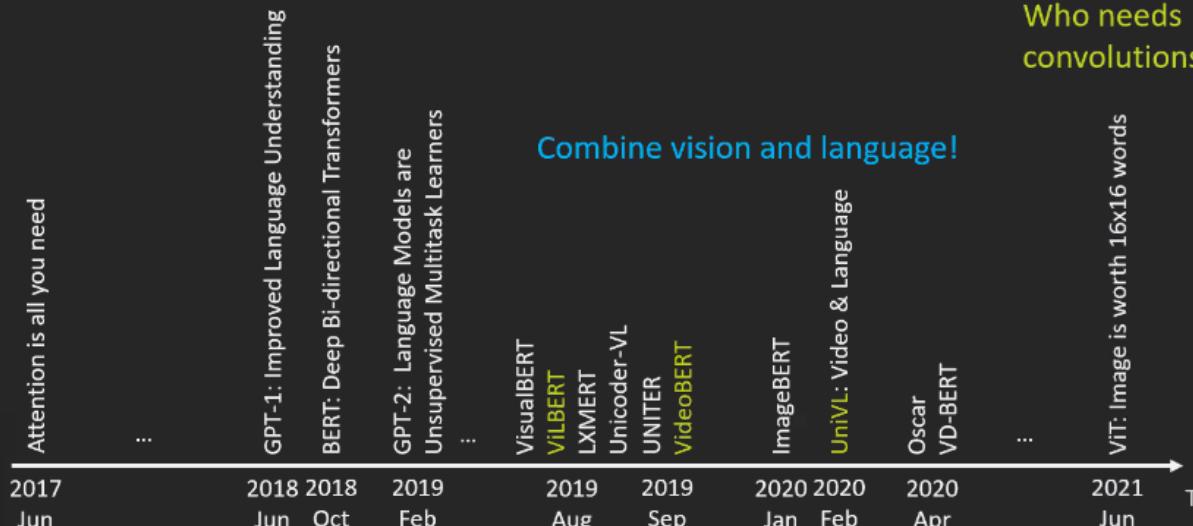
14

Transforming vision and language!



Language is “solved”!?

Who needs convolutions!?



Intro – Representations – Stories – Instructional Videos

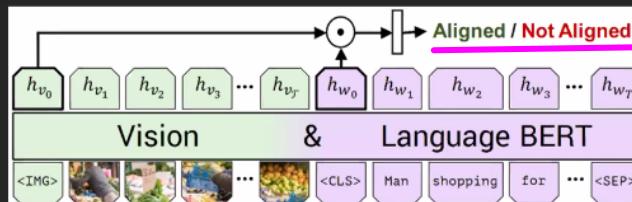
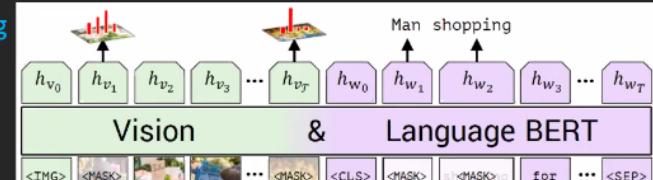
15

* To
read →

ViLBERT: Vision and Language BERT

Large-scale pretraining on image-caption pairs (Conceptual Captions)

a) Masked multi-modal learning



b) Multi-modal alignment prediction

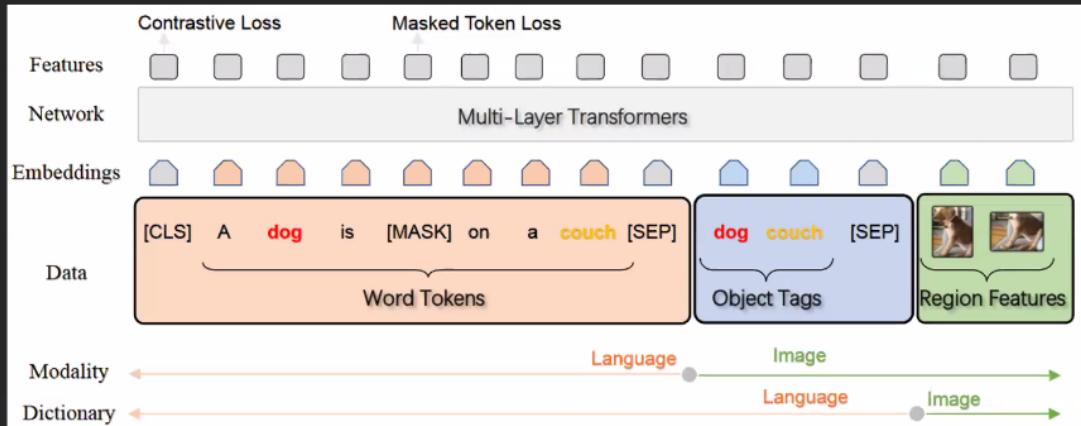
Reference: ViLBERT: Pretraining task-agnostic visio-linguistic representations for vision-and-language tasks, NeurIPS 2019

Intro – Representations – Stories – Instructional Videos

16

¥

Oscar: Image regions, Object tags, Sentence



Reference: Oscar: Object-Semantics Aligned Pretraining for Vision-Language Tasks, ECCV 2020

Intro – Representations – Stories – Instructional Videos

17

Results:

Method	Size	1K Test Set			5K Test Set		
		R@04	R@05	R@10	R@04	R@05	R@10
DVSA [14]	-	38.4	69.9	80.5	27.4	60.2	74.8
VNLBERT [7]	-	64.7	95.9	97.0	41.3	81.2	90.5
DPC [46]	-	65.0	89.8	94.9	47.1	79.9	90.0
CAMP [12]	-	72.3	94.8	98.3	58.5	87.9	95.0
SCAN [18]	-	72.7	94.8	98.4	58.8	88.4	94.8
SCG [33]	-	76.4	96.3	99.2	61.4	88.9	95.1
UNITER [1]	B	76.5	96.3	99.0	61.6	89.6	95.2
Unified-VL [19]	B	84.3	97.3	99.3	69.4	93.5	96.3
12-in-1 [24]	B	-	-	-	65.2	91.0	96.2
UNITER [5]	B	-	-	-	-	-	-
UNITER [5]	L	-	-	-	-	-	-
Oscar	B	88.4	99.1	99.8	75.7	95.2	98.3
Oscar	L	89.8	98.8	99.7	78.2	95.8	98.3

(a) Image-text retrieval

Method MAC VisualBERT LXMBERT 12-in-1 UNITER _B UNITER _L OSCAR _B OSCAR _L									
Test-dev	70.63	70.50	70.80	72.42	73.15	72.27	73.24	73.16	73.61
Test-std	70.92	70.83	71.00	72.54	-	72.46	73.40	73.44	73.82

(b) VQA

Method MAC VisualBERT LXMBERT 12-in-1 UNITER _B UNITER _L OSCAR _B OSCAR _L									
Dev	50.8	67.40	74.90	-	77.14	78.40	78.07	79.12	-
Test-P	51.4	67.00	74.50	78.87	77.87	79.50	78.36	80.37	-

(c) NLVR2

Method cross-entropy optimization CIDEr optimization									
		M	C	S	M	C	S		

Method cross-entropy optimization CIDEr optimization									
		M	C	S	M	C	S		

Method in-domain CIDEr SPICE near-domain CIDEr SPICE out-of-domain CIDEr SPICE overall CIDEr SPICE									
		M	C	S	M	C	S		

Method in-domain CIDEr SPICE near-domain CIDEr SPICE out-of-domain CIDEr SPICE overall CIDEr SPICE									
		M	C	S	M	C	S		

Method in-domain CIDEr SPICE near-domain CIDEr SPICE out-of-domain CIDEr SPICE overall CIDEr SPICE									
		M	C	S	M	C	S		

Method in-domain CIDEr SPICE near-domain CIDEr SPICE out-of-domain CIDEr SPICE overall CIDEr SPICE									
		M	C	S	M	C	S		

Method in-domain CIDEr SPICE near-domain CIDEr SPICE out-of-domain CIDEr SPICE overall CIDEr SPICE									
		M	C	S	M	C	S		

Method in-domain CIDEr SPICE near-domain CIDEr SPICE out-of-domain CIDEr SPICE overall CIDEr SPICE									
		M	C	S	M	C	S		

Method in-domain CIDEr SPICE near-domain CIDEr SPICE out-of-domain CIDEr SPICE overall CIDEr SPICE									
		M	C	S	M	C	S		

Method in-domain CIDEr SPICE near-domain CIDEr SPICE out-of-domain CIDEr SPICE overall CIDEr SPICE									
		M	C	S	M	C	S		

Method in-domain CIDEr SPICE near-domain CIDEr SPICE out-of-domain CIDEr SPICE overall CIDEr SPICE									
		M	C	S	M	C	S		

Method in-domain CIDEr SPICE near-domain CIDEr SPICE out-of-domain CIDEr SPICE overall CIDEr SPICE									
		M	C	S	M	C	S		

Method in-domain CIDEr SPICE near-domain CIDEr SPICE out-of-domain CIDEr SPICE overall CIDEr SPICE									
		M	C	S	M	C	S		

Method in-domain CIDEr SPICE near-domain CIDEr SPICE out-of-domain CIDEr SPICE overall CIDEr SPICE									
		M	C	S	M	C	S		

Method in-domain CIDEr SPICE near-domain CIDEr SPICE out-of-domain CIDEr SPICE overall CIDEr SPICE									
		M	C	S	M	C	S		

Method in-domain CIDEr SPICE near-domain CIDEr SPICE out-of-domain CIDEr SPICE overall CIDEr SPICE									
		M	C	S	M	C	S		

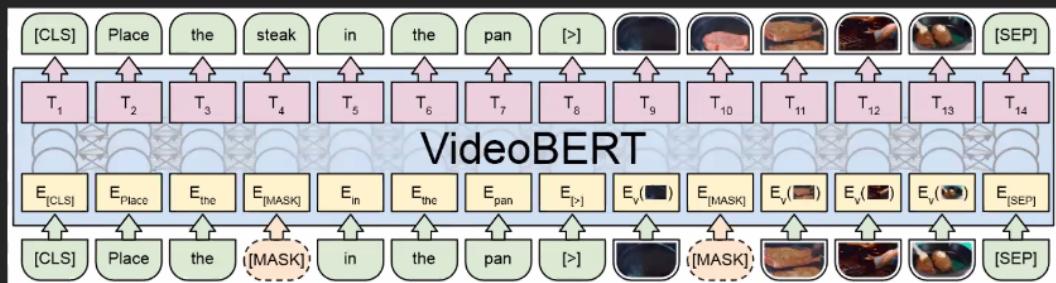
Method in-domain CIDEr SPICE near-domain CIDEr SPICE out-of-domain CIDEr SPICE overall CIDEr SPICE									
		M	C	S	M	C	S		

Method in-domain CIDEr SPICE near-domain CIDE

VideoBERT: Video and language, not different!



*Video &
Text
Learning
{Images → frames}

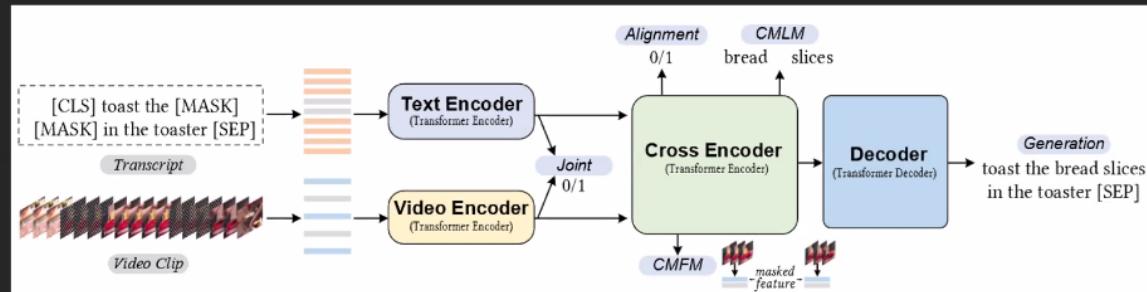


Reference: VideoBERT: A Joint Model for Video and Language Representation Learning, ICCV 2019

Intro – Representations – Stories – Instructional Videos

19

UniVL



Reference: UniVL: A Unified Video and Language Pre-Training Model for Multimodal Understanding and Generation, arXiv 2020.02

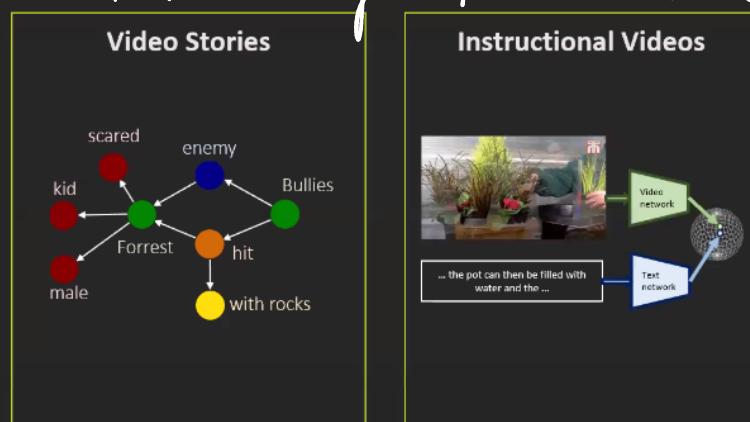
Intro – Representations – Stories – Instructional Videos

20

Two main streams in this talk

Video: Forest Gump chased by his friends with rocks.

Preparing
Graphs to
understand the
video



[Not the sequence for now]

- Overview -

Stories naturally engage vision and language

Human
Interpretation →
of Stories



Drawing credit: Ananya Tapaswi, age 6

High in the mountains,
lived a fire breathing
dragon. Terrorizing the
people of a nearby
village had become its
daily ritual. Fed up and
frightened, the villagers
finally go to their
queen and ask for help.

Stories reflect human behavior



Panchatantra

Features stories about:

- Winning friends
- Losing friends
- Rash actions leading to losses
- How to face difficult situations
- Strategies of war and peace

Shakespeare

Features stories about:

- romance (Romeo & Juliet)
- manipulation and deceit (Macbeth)
- humor (A Midsummer Night's Dream)
- envy (The Merchant of Venice)

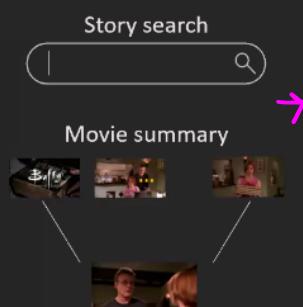


Story understanding

The ability of a machine to “watch” or “read” a story and perform human-like tasks



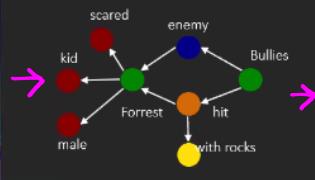
Search & Summarize



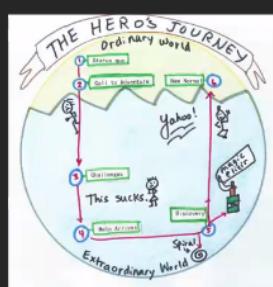
Answer Questions



Learn Human Behavior



Create new stories!



Learning from Movie texts

Transcripts and Subtitles with videos

Leonard: At least I didn't have to invent twenty-six dimensions just to make the math come out.
Leonard stands up and looks at Sheldon.
Sheldon: I didn't invent them, they're there.
Leonard: In what universe?.

00:07:00,733 --> 00:07:04,612
At least I didn't have to invent 26 dimensions to make the math come out.
00:07:04,773 --> 00:07:07,412
- I didn't invent them. They're there.
- In what universe?

Hello! My name is...
Buffy" – Automatic
Naming of Characters
in TV Video,
BMVC 2006



Learning realistic
human actions
from movies,
CVPR 2008

How to learn from multiple depictions



Plot Synopsis

Noting that the dire wolf is the sigil of the Stark family and there are as many pups as the Stark children, they take the pups in as companions.

Back at Winterfell, Catelyn informs her husband of a letter announcing the death of Lord Arryn, Eddard's old mentor and Catelyn's brother-in-law.

An additional message reports that the king himself is coming to Winterfell.



Film / TV episode



Book



"You did not come here to tell me crib tales. I know how little you like this place. What is it, my lady?"

Catelyn took her husband's hand. "There was grievous news today, my lord. I did not wish to trouble you until you had cleansed yourself." There was no way to soften the blow, so she told him straight. "I am so sorry, my love. Jon Arryn is dead."

His eyes found hers, and she could see how hard it took him, as she had known it would. In his youth, ...

Search by
Scene descriptions →

Glory sucks Taras mind

You searched for: Glory sucks Taras mind

time: 1661.284

Play

Score	Video	Line No(s.)	Matched Sentences
15.32	Buffy S05E19 ► 0:27:24 - 0:27:54 ► 0:27:13 - 0:27:51	15	Protecting Dawn, Tara refuses, and Glory drains Tara's mind of sanity
13.82	Buffy S05E19 ► 0:25:50 - 0:27:54	14, 15	Glory discovers that Tara isn't the Key, and offers to let her go if she reveals the key's identity. Protecting Dawn, Tara refuses, and

Search based on story plots

The screenshot shows a search interface with a search bar containing "Glory sucks Tara's mind". Below the search bar, it says "You searched for: Glory sucks Tara's mind". A video thumbnail is displayed, showing a scene from a TV show where a character is shouting. Below the video, there is a timeline with markers at 27.47 and 42.44. The text "Time: 1667 560" and a "Hide" button are also visible. At the bottom, a table lists two search results:

Score	Video	Line No(s.)	Matched Sentences
15.32	Buffy S05E19 ► 0:27:24 - 0:27:54 ► 0:27:19 - 0:27:51	15	Protecting Dawn, Tara refuses, and Glory drains Tara's mind of sanity
13.82	Buffy S05E19 ► 0:25:59 - 0:27:54	14, 15	Glory discovers that Tara isn't the Key, and offers to let her go if she reveals the key's identity. Protecting Dawn, Tara refuses, and

Buttons for "Play!" are shown next to each result.

62 queries

85% found in plot

15% not localized in plot

53 queries

40 have overlap

13 don't overlap

[Tapaswi, IJMIR 2015]

Intro – Representations – Stories – Instructional Videos

31

Generate narrative captions

Plot synopses

1 sentence, video ~60+ s



Buffy performs the ritual in her bedroom then walks around her house looking for anything unusual.

[Tapaswi, IJMIR 2015]

Books

1 paragraph, video ~20 s



From an inside pocket of his black overcoat he pulled a slightly squashed box. Harry opened it with trembling fingers. Inside was a large, sticky chocolate cake with Happy Birthday Harry written on it in green icing.

[Tapaswi, CVPR 2015]

Intro – Representations – Stories – Instructional Videos

32

Characters are important!

33

MovieQA: Story question-answering

The Matrix has you
00:25:52 -> 00:25:57 Welcome, Neo. As you no doubt have guessed... I am Morpheus
00:40:42 -> 00:40:47 It exists now only as part of a neural-interactive simulation that we call the Matrix.
01:04:08 -> 01:04:09 ... you know what I realize?
Ignorance is bliss.
02:08:38 -> 02:08:39 Where we go from there is a choice I leave to you

Movie:

- 200,000 frames
- 2,000 shots
- 1,000 dialogs
- Long temporal dependencies
- Actions, interactions, emotions, intent

[Tapaswi, CVPR 2016]

Intro – Representations – Stories – Instructional Videos

34

Identifying characters is crucial

*
Character
Identification
Crucial in
movies.

- Q. Who makes Indy return the crucifix after escaping from the grave robbers?
- A1. The local sheriff
 - A2. Coronado
 - A3. No one, he keeps it
 - A4. The Boy Scout troop
 - A5. The grave robbers



Intro – Representations – Stories – Instructional Videos

35

How do we know characters' names?



- ✓ • Dialog between characters
“Hi, I’m Sheldon”
“Raj, can you pass me the water?”
- ✓ • Use weak annotations from dialog
- ✓ • Multiple instance learning

[Haurilet, WACV 2016]

Intro – Representations – Stories – Instructional Videos

36

More than facial appearance



- Identity is more than just faces
- Learn **clothing** on-the-fly
- Speech patterns and **speaking styles**
- Energy based model to capture interactions across modalities

[Tapaswi, CVPR 2012]

Intro – Representations – Stories – Instructional Videos

37

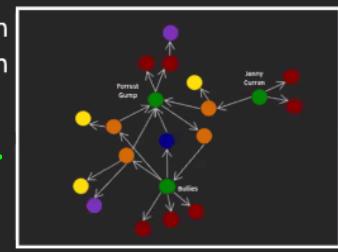


Video clips of social situations



time

Situation Graph



To understand
Characters
behaviour

Intro – Representations – Stories – Instructional Videos

40

What are main aspects of human behaviour ?

Situations	argument	wedding	introduction	business meeting
	work talk	fight	phone call	intimacy
Interactions	asks	orders	listens	leaves
	hits	explains	kisses	advises
Relationships	neighbor	friend	lover	parent
	stranger	colleague	spouse	customer
Attributes	female	adult	serious	nervous
	doctor	worried	happy	upset
			quiet	confused

Intro – Representations – Stories – Instructional Videos

41

Tags
for
human
behaviour

What are MovieGraphs?

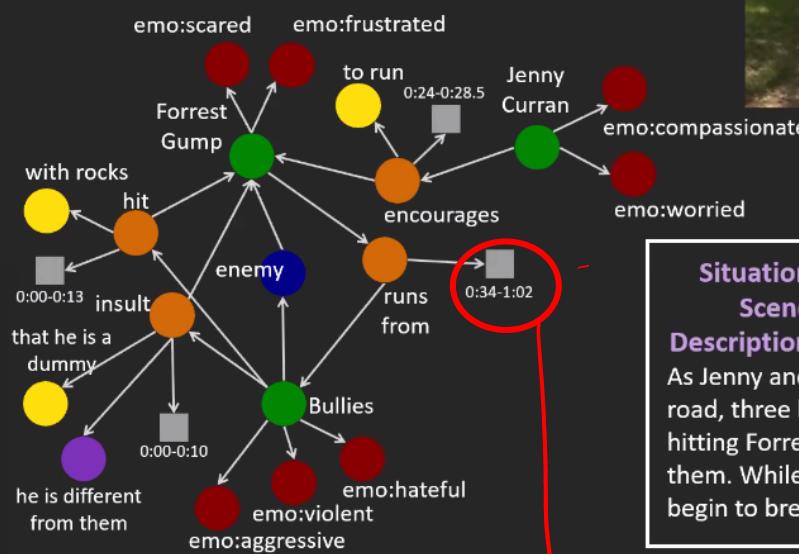
Film Credit:
Forrest Gump



Intro – Representations – Stories – Instructional Videos

42

What are MovieGraphs?



Situation: Bullying
Scene: Field road
Description:

As Jenny and Forrest are walking down the road, three boys come along and start hitting Forrest. Jenny urges him to run from them. While Forrest runs, his leg braces begin to break apart.

Intro – Representations – Stories – Instructional Videos

42

→ Adding time stamps → sequence of events

MovieGraphs Dataset



Exploring Social Common Sense

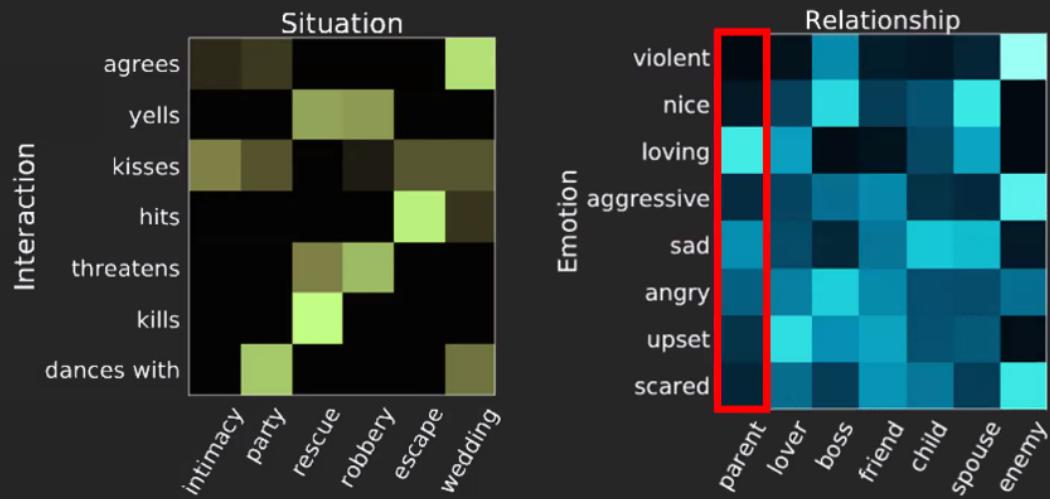
Social situations can be tricky!

**"I'm sorry" and "my bad"
mean the same thing...**

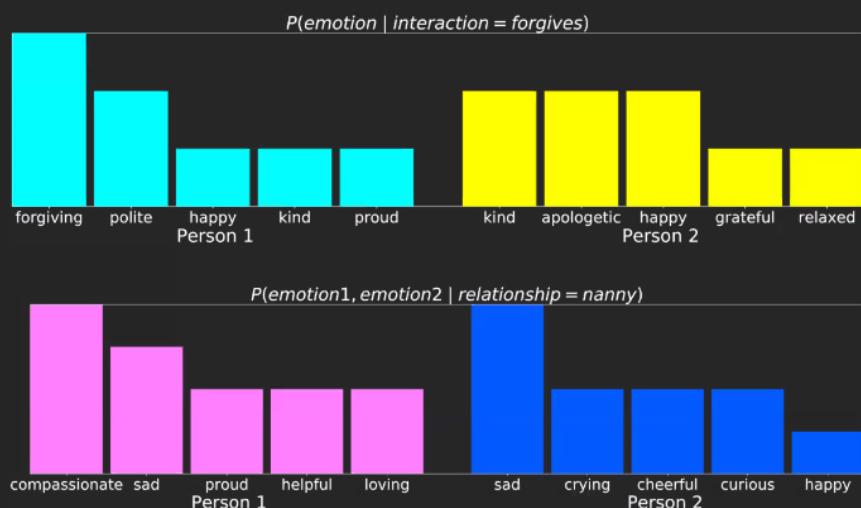


**Unless you are at
a funeral.**

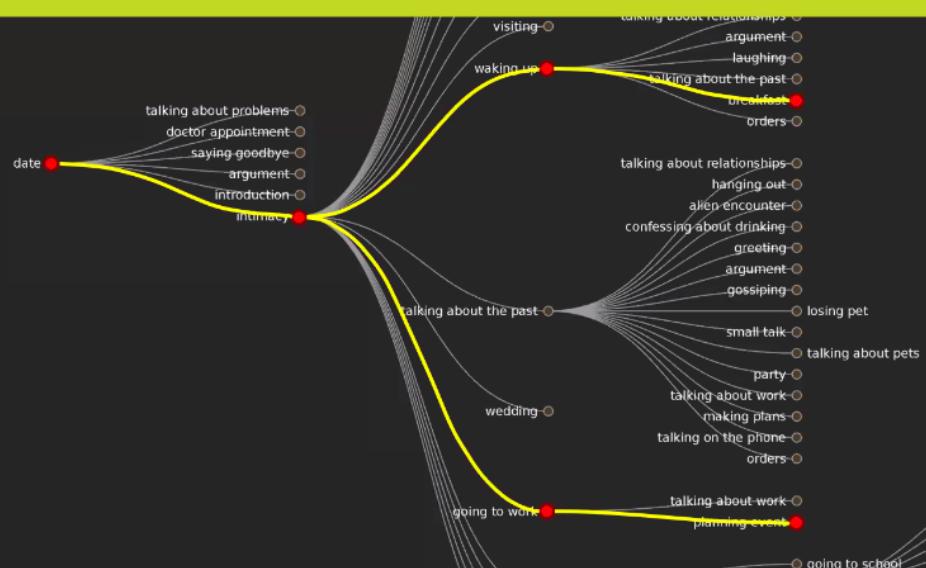
What happens when ...



Conditional Attributes



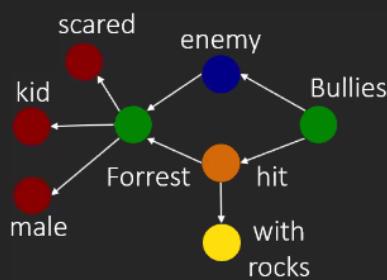
Situation Flow



Graph based Video Retrieval

49

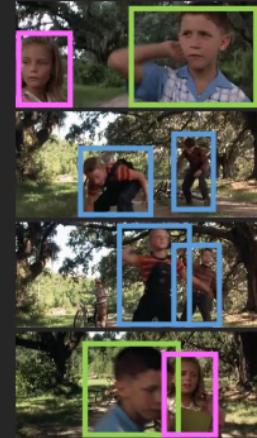
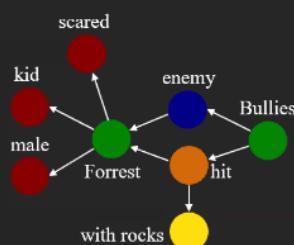
Graph-based Video Retrieval



Intro – Representations – Stories – Instructional Videos

Scoring similarity: graph query and video

$$F(G, M, z) = \sum_c \phi(g_c, m_c, z)$$



Intro – Representations – Stories – Instructional Videos

Query Based Retrieval of Scenes.

Graph-based clip retrieval results - I

Query
Scene: Garden
Situation: Wedding

Rank 1



Jerry Maguire, 1996

Rank 2



Meet the Parents, 2000

Query
Jacob
Emily
Cal
David

Crazy, Stupid, Love., 2011

Rank 1



Jacob, Emily, Cal, David

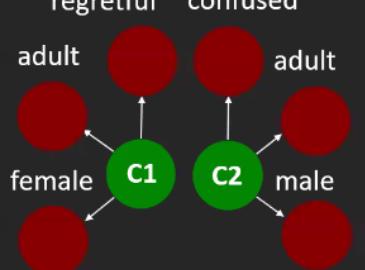
Rank 2



Jacob, Cal

Intro – Representations – Stories – Instructional Videos

Graph-based clip retrieval results - II

Query
Situation: Relationship talk
Scene: Church
regretful confused
adult adult
female male


Rank 1



The Firm, 1993

situation: argument
scene: banquet hall
C1: adult, female,
heartbroken, angry
C2: adult, male,
honest, remorseful

Rank 2



Four Weddings and a
Funeral, 1994

situation: relationship talk
scene: church
C1: adult, female,
regretful
C2: adult, male,
confused

Intro – Representations – Stories – Instructional Videos

53

Interactions and Relationships

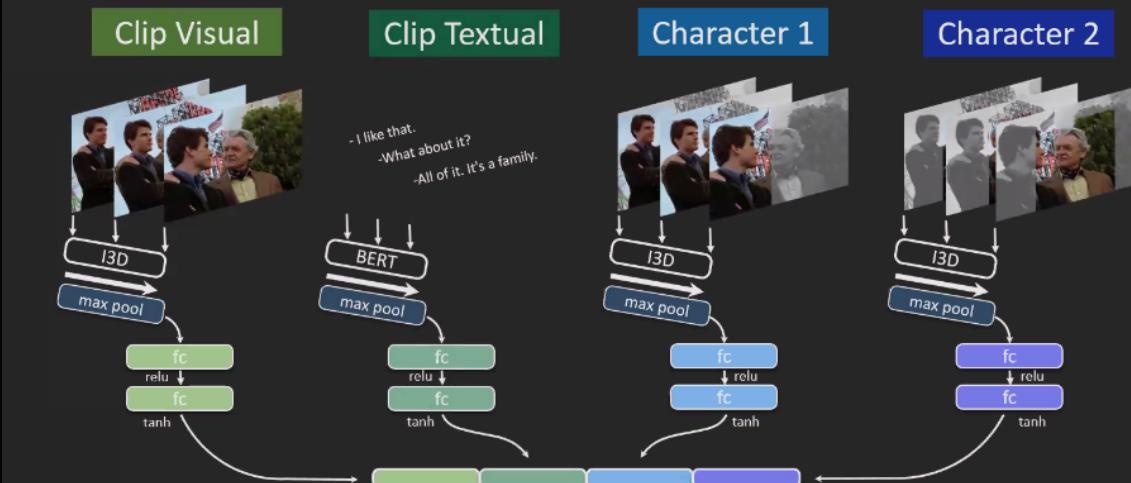
54

A joint study of *interactions* and *relationships*

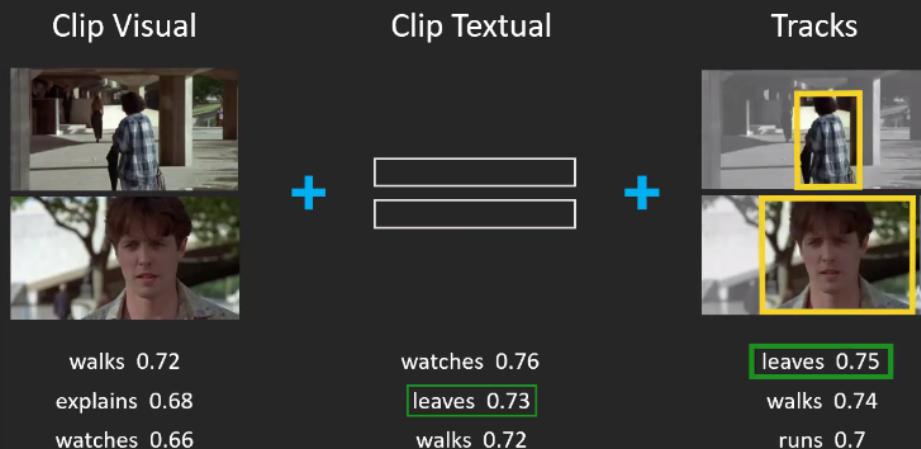


[Kukleva, CVPR 2020]

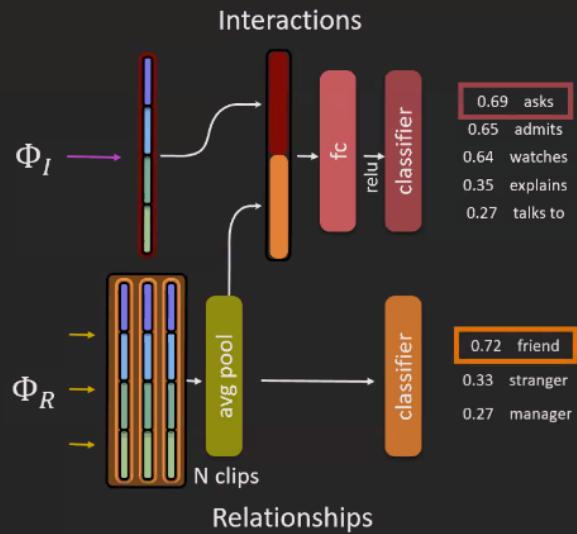
Multimodal embedding



Interaction prediction with multiple modalities



Interactions, relationships, and multi-task learning



Joint prediction

Model	Interaction	Relationship
Random	1.0	6.7
Interaction	26.1	-
Relationship	-	26.8
Joint	26.3 (+0.2)	28.1 (+1.3)

Interactions

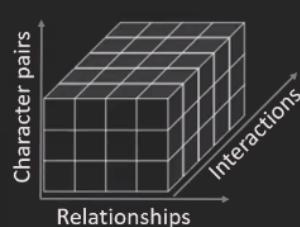
Hugs (siblings, child, lover)	+17%
Introduces (colleague, friend, stranger)	+14%
Runs (enemy, lover)	+12%
Talks to (all relationships)	-1%
Informs (many relationships)	-5%

Relationships

Sibling (hugs, informs, explains, watches)	+10%
Acquaintance (explains, greets)	+8%
Lover (suggests, kisses)	+7%
Parent (hugs, watches, suggests, explains)	-5%
Manager (talks to, suggests)	-8%

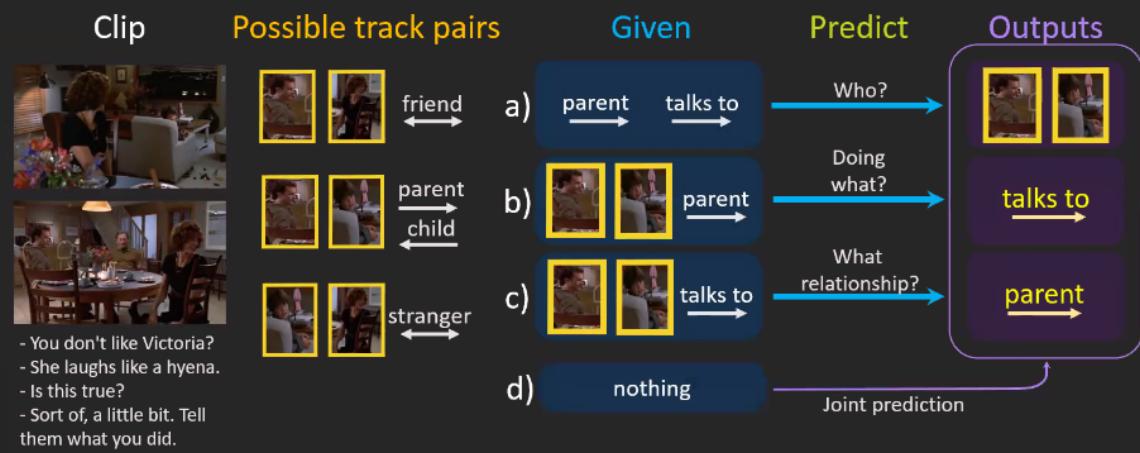
Which pair of characters are involved?

- Training with weak clip-level labels
- Unknown character pair is treated as a latent variable
- A max-margin loss with two tweaks is used to train models:
 - Later training epochs use hard negative sampling
 - Estimate \hat{p} is sampled from a probability distribution rather than argmax
- At inference, we obtain a 3D tensor and compute argmax



Task	Given	Predict	Slice
Who?	interaction relationship	character	
What relationship?	character interaction	relationship	
Doing what?	character relationship	interaction	

Who, doing what, and what relationship?

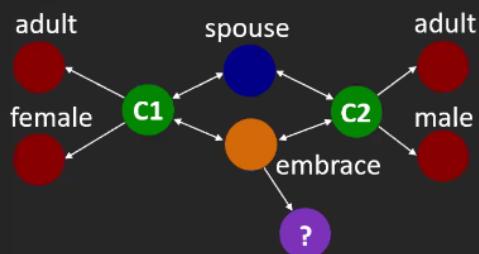


[Kukleva, CVPR 2020]

Reason prediction

Scene: House

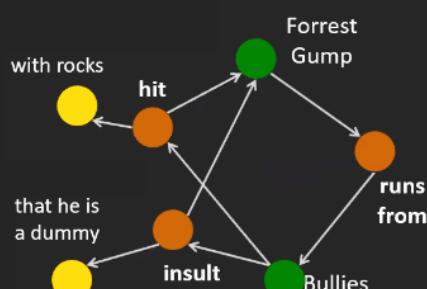
Situation: Talk about adoption



PRED: he is scared

GT: glad they'll have a baby

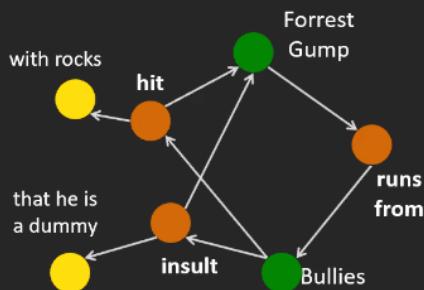
Interaction ordering



Unordered interactions

- | | |
|---------|---------------------------|
| Forrest | runs from |
| Bullies | hit with rocks |
| Bullies | insult that he is a dummy |

Interaction ordering



Ordered interactions

Bullies insult that he is a dummy

Bullies hit with rocks

Forrest runs from

t
i
m
e
↓

Revival of TV / Movie Understanding

Large Scale Movie Description Challenge (LSMDC)



His brow furrowed,
SOMEONE looks down
at the ground.

SOMEONE eyes him
angrily, her jaw
clenched.

SOMEONE heads off.

SOMEONE folds her
arms.

SOMEONE approaches
SOMEONE, who leans against
the wall of the house.

Just
identified →
all human's
as person

Large Scale Movie Description Challenge (LSMDC)



His brow furrowed,
[PERSON1] looks down
at the ground.

[PERSON2] eyes him
angrily, her jaw
clenched.

[PERSON1] heads off.

[PERSON2] folds her
arms.

[PERSON1] approaches
[PERSON3], who leans against
the wall of the house.

LSMDC, 2015, <https://sites.google.com/site/describingmovies/>

Intro – Representations – Stories – Instructional Videos

66

Condensed Movies



Condensed Movies, ACCV 2020

<https://www.robots.ox.ac.uk/~vgg/research/condensed-movies/>

Intro – Representations – Stories – Instructional Videos

67

MovieNet

*
Japan
→
Researcher
Artist.

The MovieNet interface displays the following sections for the movie *Titanic*:

- Photo:** Includes a poster for *Titanic* and a photo of Leonardo DiCaprio and Kate Winslet.
- Movie:** A horizontal strip of several video frames from the movie.
- Trailer:** A horizontal strip of video frames from the movie trailer.
- Subtitle:** A table showing subtitle entries with timestamps and text:

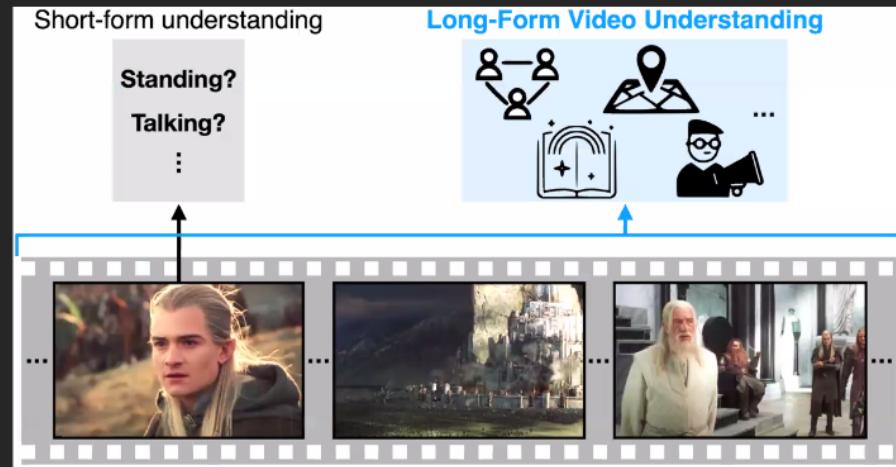
02:13:31,355 -> 02:13:34,028	You're so stupid! Why did you do that?
02:13:34,435 -> 02:13:36,107	You're so stupid, Rose.
02:13:37,835 -> 02:13:39,747	Why did you do that? Why?
02:13:40,316 -> 02:13:42,193	You jump, I jump, right?
- Meta Data:** Includes the title, runtime, genre, rating, director, cast, and storyline.
- Wiki Plot:** A brief summary of the plot.
- Synopsis:** A detailed synopsis of the scene where Rose boards the boat.
- Script:** The full dialogue script for the scene.

MovieNet, ECCV 2020, <http://movienet.site/>

Intro – Representations – Stories – Instructional Videos

68

Long-form Video Understanding

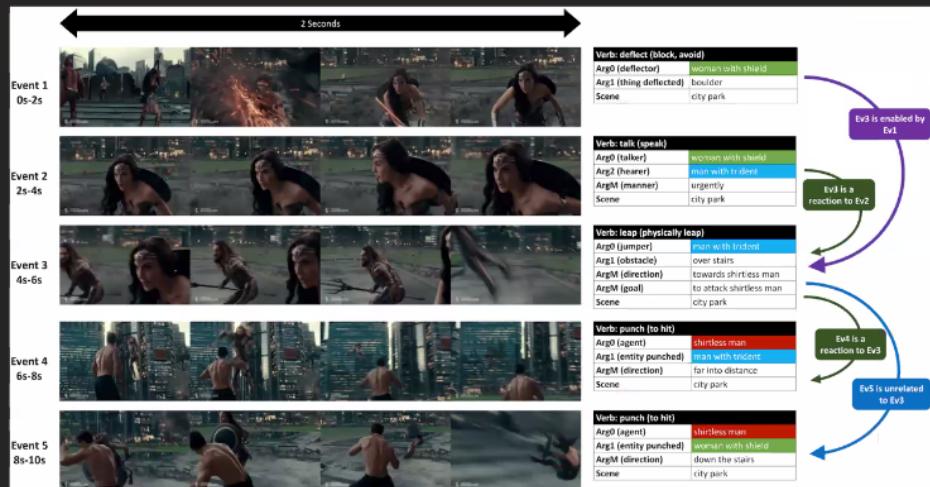


LVu, CVPR 2021, <https://chaoyuan.org/lvu/>

Intro – Representations – Stories – Instructional Videos

69

VidSitu: Video Situation Recognition



VidSitu, CVPR 2021, <https://vidsitu.org/>

Intro – Representations – Stories – Instructional Videos

70

VALUE challenge: Retrieval, QA, Captioning



Multi-channel Video
With both Video Frames and Subtitle/ASR



Diverse Video Domain
Diverse video content from YouTube, TV Episodes and Movie Clips



Various Datasets over Representative Tasks
11 datasets over 3 tasks: Retrieval, Question Answering and Captioning.



Leaderboard!
To track the advances in Video-and-Language research.

Attend Mohit Bansal's talk tomorrow!

VALUE Benchmark, <https://value-benchmark.github.io/>

Intro – Representations – Stories – Instructional Videos

71

Extracting
short freely
available
clips helps

Representation learning from Instructional videos



- Overview -

Can we learn without manual annotations?

Key idea 1
a **wide variety**
of instructional videos are available
at **large scale** on the Internet!

Key idea 2
Free supervision
using ASR output

Intro – Representations – Stories – Instructional Videos

WikiHow: filtering to visual instructional content

The screenshot shows the WikiHow homepage with a search bar containing "How to Change Your Whole Personality". Below the search bar, several article cards are displayed:

- How to Be**: Co-authored by Emily May, Updated: September 5, 2019. Description: There's nothing like the own kitchen. Baking a c ingredients, mixing ther remembering to take th burns. Read on to learn pound cake, chocolate.
- How to C**: Co-authored by Andrew C, Updated: August 27, 2019. Description: Have you ever been s life? Do you want to b to ask for help? Fortu simple task, provided little effort.
- How to Grow Vegetables in Small**: Co-authored by Andrew Carrberry, Updated: March 29, 2019. Description: Even the smallest of vegetable gardens can yield big returns. If you are limited on space but still want to enjoy fresh vegetables, you can use a variety of techniques to ensure a plentiful harvest all season long.
- How to Change Your Whole Personality**: Co-authored by wikiHow Staff, Updated: August 8, 2019. Description: Personality is a collection of patterns — thought, behavior, and feeling — that make up who you are. And guess what? Patterns can change. It'll take work, but if you're truly devoted to this idea, anything can happen. Remember, though, that your old personality will likely shine through regularly as our beliefs and thinking is shaped by our life experiences.

Below the cards, a summary states: "Scraped 130K tasks, kept 23K with potential visual cues".

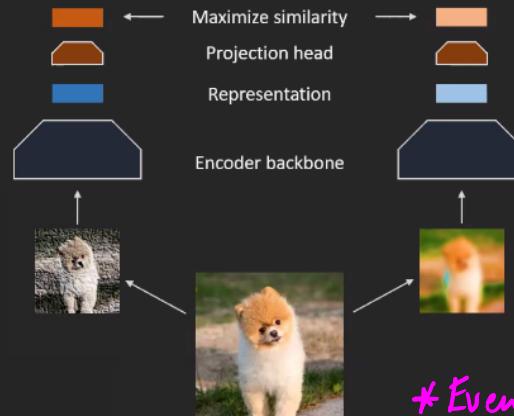
Intro – Representations – Stories – Instructional Videos



- Videos of people performing over 23k visual tasks
- 1.23 million Youtube videos (15 years!)
- Contains 130 millions of ~4s clips with narration
- Larger than any annotated video dataset

[Miech, ICCV 2019]

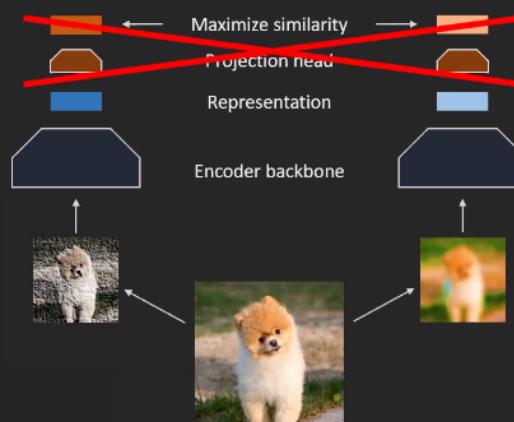
SimCLR: modern contrastive learning



*Even if labels are not known at the first place

Reference: A Simple Framework for Contrastive Learning of Visual Representations, ICML 2020

SimCLR: modern contrastive learning

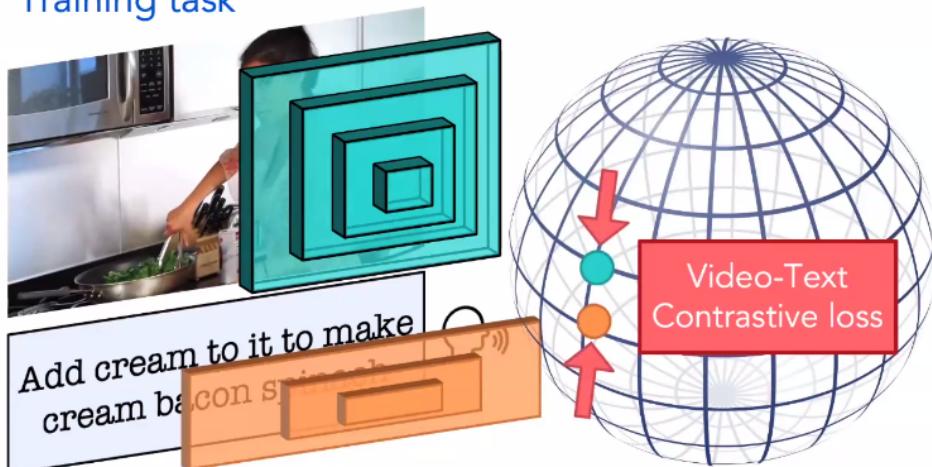


Reference: A Simple Framework for Contrastive Learning of Visual Representations, ICML 2020

Video-text contrastive loss

Video Credit: End-to-End Learning of Visual Representations from Uncurated Instructional Videos, CVPR 2020

Training task



But the clips are not aligned ...

Video Credit: End-to-End Learning of Visual Representations from Uncurated Instructional Videos, CVPR 2020

But what
about in
practice ?

But the clips are not aligned ...

Video Credit: End-to-End Learning of Visual Representations from Uncurated Instructional Videos, CVPR 2020

Speech
Recognition



you can add cilantro
basil they give

But the clips are not aligned ...

Video Credit: End-to-End Learning of Visual Representations from Uncurated Instructional Videos, CVPR 2020

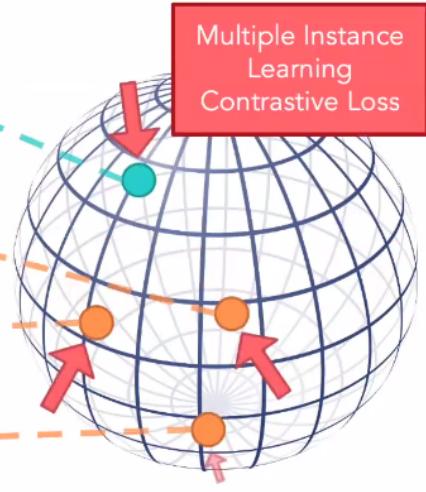


But the clips are not aligned ...

Video Credit: End-to-End Learning of Visual Representations from Uncurated Instructional Videos, CVPR 2020



- Spinach what's the name
- Give it a couple more tosses
- Fresh herbs maybe Some oregano



Multiple Instance Learning – Noise Contrastive Estimation

a) Standard Contrastive Loss

$$L_{i,j} = -\log \frac{\exp(\text{sim}(z_i, z_j))}{\sum_{k \neq i} \exp(\text{sim}(z_i, z_k))} \quad * \text{Typo}$$

b) MIL-NCE

$$L_{i,P_i} = -\log \frac{\sum_{j \in P_i} \exp(\text{sim}(z_i, z_j))}{\sum_{j \in P_i} \exp(\text{sim}(z_i, z_j)) + \sum_{k \in N_i} \exp(\text{sim}(z_i, z_k))}$$

ActBERT

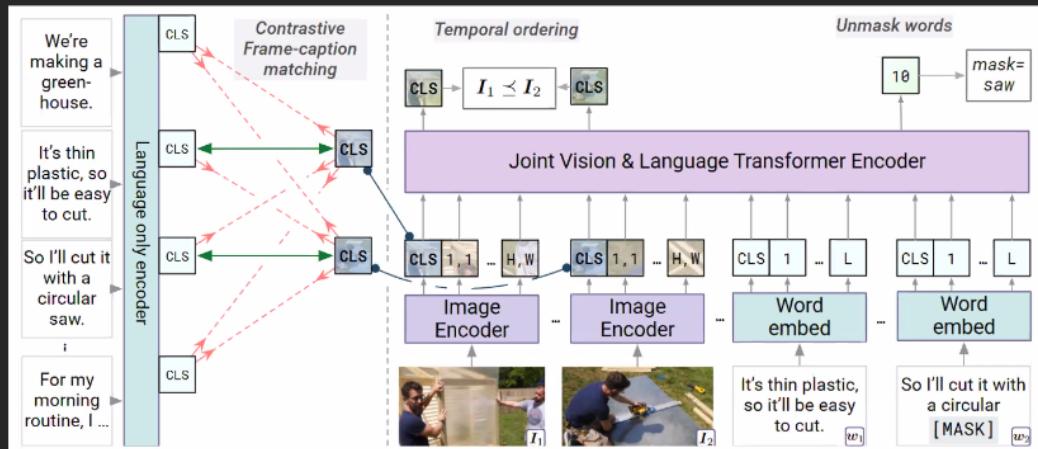


Reference: ActBERT: Learning Global-Local Video-Text Representations, CVPR 2020

Intro – Representations – Stories – **Instructional Videos**

81

MERLOT: Multimodal knowledge models



Reference: MERLOT: Multimodal Neural Script Models, arXiv 2021.06

Intro – Representations – Stories – **Instructional Videos**

82

Thanks to my students, collaborators, advisors!



83

If any of this is interesting for you ...

Reach out!

I'm looking for motivated students to work
on video-language problems

<https://makarandtapaswi.github.io/students/>

84



Thank you!

Makarand Tapaswi

✉ makarand.tapaswi@iiit.ac.in

🌐 <https://makarandtapaswi.github.io/>