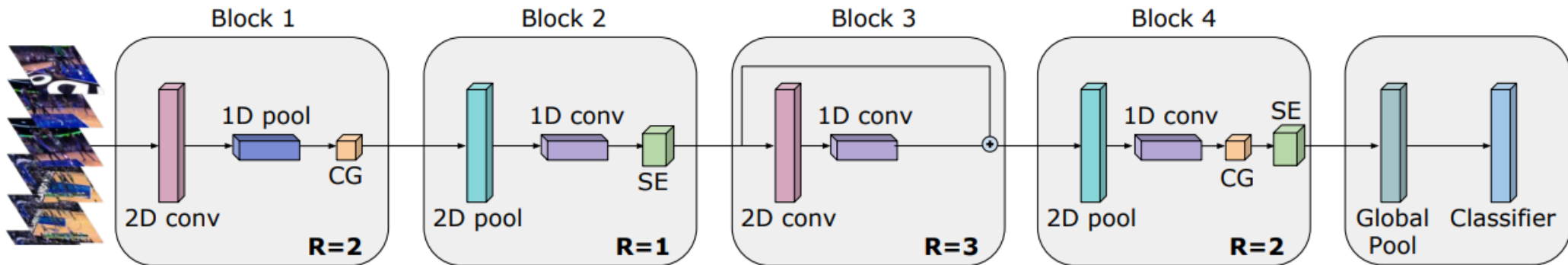
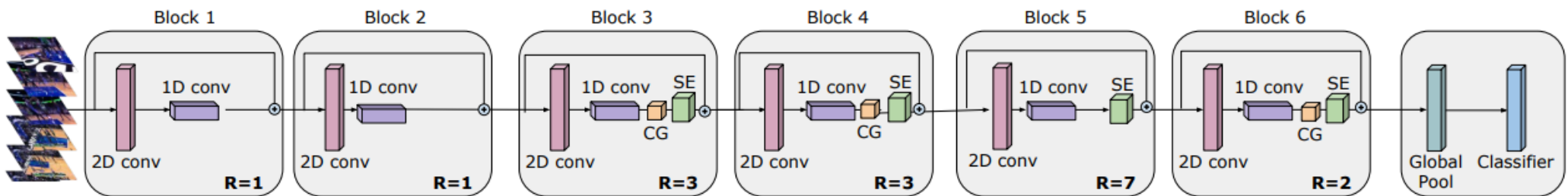


# Tiny Video Networks on UCF50 dataset

Input Video



Input Video

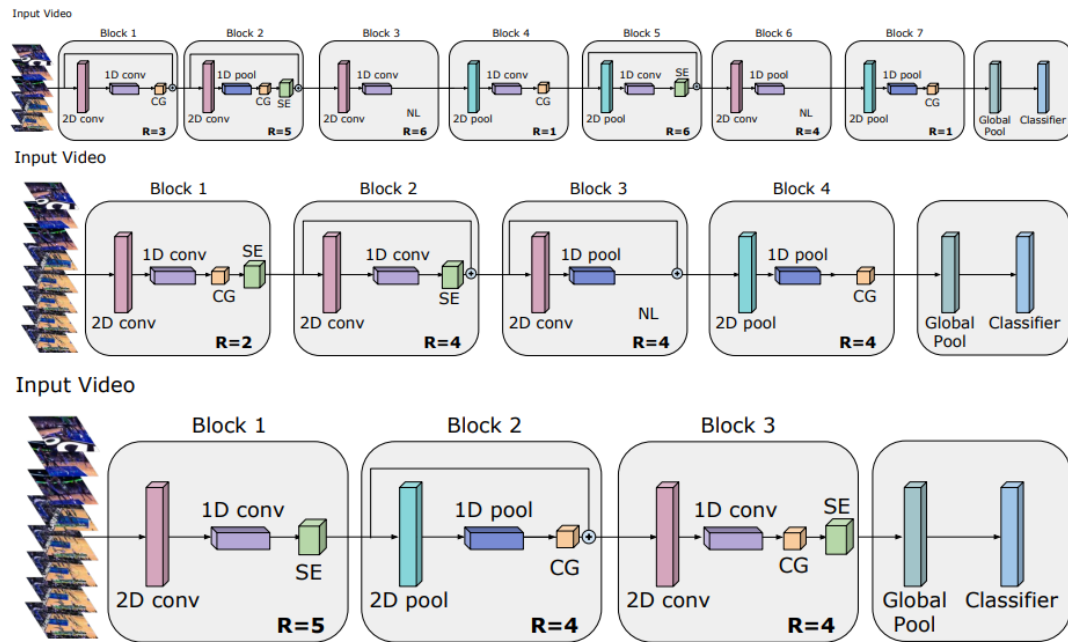


Выполнила Чувиляева Ирина

# Постановка задачи

- Воспроизвести одну из сетей, описанный в статье Tiny Video Networks, AJ Piergiovanni Anelia Angelova Michael Ryoo
- Обучить её на датасете UCF50
- Сравнить с известными результатами

# О статье



3: Example Tiny Video Networks found using architecture evolution showing several blocks with different configurations. A Tiny Video Net has multiple blocks, each repeated  $R$  times. Each block has a different configuration with spatial and temporal convolution, pooling, non-local layers, context gating and squeeze-excitation layers. It can select the image resolution and frame rate. From top to bottom: TVN-2, TVN-3, TVN-4. TVN-1 is shown in Figure 1.

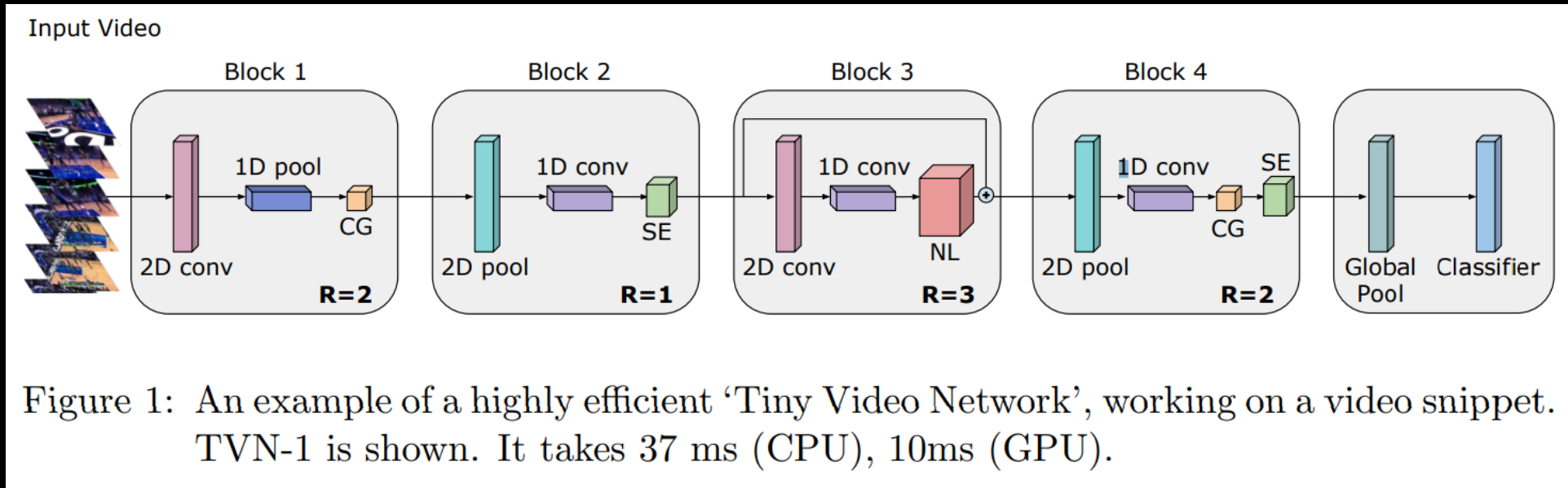
Method	Runtime CPU(ms)	Runtime GPU(ms)	GFlops	Acc. (%)
ResNet-18	2120	105	38	21.1
ResNet-34	2256	110	50	24.2
ResNet-50	3022	125	124	28.1
ResNet-101	3750	140	245	30.2
TSN [47]	-	-	-	24.1
2DResNet50 [26]	-	-	-	27.1
bLVNet-TAM [7]	-	-	-	<b>31.4</b>
TVN-1	37	10	13	23.1
TVN-2	65	13	17	24.2
TVN-3	85	16	69	25.9
TVN-4	402	19	106	28.0
TVN-5	86	16	52	<b>29.8</b>
TVN-6	142	18	93	<b>30.7</b>

Table 3. Results on the MiT dataset comparing different Tiny Networks to baselines and state-of-the-art (which are all RGB-only). TVN models perform competitively and are also much faster. No runtime was reported in prior work.

# О датасете UCF50



# О реализации



- Архитектура TVN1
- Входной размер 200x200
- Нет аугментации при обучении
- Нет предобучения
- Соотношение train/test: 85/15
- Берём каждый 10ый кадр из видео
- Всего берём 8 кадров (из 80)
- Если кадров не хватает – повторяем последний

The screenshot shows the JupyterLab web application running in a browser at localhost:8888. The notebook is named "train\_tvn1\_framed". The top toolbar includes icons for file operations, running code, and switching kernels. Below the toolbar, the notebook output displays the execution of a training script. It shows progress bars for epochs 6 through 11, each reaching 100% completion. Each epoch's output includes the mean loss as a tensor value and the accuracy as a float. The final line shows the progress bar at 35%.

```
100% |████████████████████████████████████████| 1420/1420 [14:32<00:00, 1.63it/s]
epoch 6 :
mean loss = tensor(3.9050)

100% |████████████████████████████████████████| 251/251 [01:44<00:00, 2.41it/s]
mean val loss = tensor(3.9034)
accuracy = 0.0037387836490528413

100% |████████████████████████████████████████| 1420/1420 [14:29<00:00, 1.63it/s]
epoch 7 :
mean loss = tensor(3.9050)

100% |████████████████████████████████████████| 1420/1420 [14:33<00:00, 1.63it/s]
epoch 8 :
mean loss = tensor(3.9050)

100% |████████████████████████████████████████| 251/251 [01:46<00:00, 2.36it/s]
mean val loss = tensor(3.9029)
accuracy = 0.0037387836490528413

100% |████████████████████████████████████████| 1420/1420 [14:37<00:00, 1.62it/s]
epoch 9 :
mean loss = tensor(3.9045)

100% |████████████████████████████████████████| 1420/1420 [14:28<00:00, 1.63it/s]
epoch 10 :
mean loss = tensor(3.9049)

100% |████████████████████████████████████████| 251/251 [01:44<00:00, 2.39it/s]
mean val loss = tensor(3.9056)
accuracy = 0.0037387836490528413

100% |████████████████████████████████████████| 1420/1420 [14:31<00:00, 1.63it/s]
epoch 11 :
mean loss = tensor(3.9048)

35% |██████████████████
```



# Сравнение с другими моделями

Performance	Experimental Setup	Paper
76.90%	Leave One Group Out Cross-validation (25 cross-validations)	Reddy and Shah. (MVAP), 2012
57.90%	5-fold group-wise cross-validation	Sadanand and Corso. (CVPR), 2012
76.40%*	Video Wise Cross-validation (*Since videos belonging to a group are obtained from a single long video, similar videos can end up in both training and testing in "video-wise cross-validation" leading to high performance)	Sadanand and Corso. (CVPR), 2012
81.03%*	2/3 training and 1/3 testing for each class (*From the details given in the paper, we are not sure if videos belonging to the same group are kept separate in training and testing sets and the paper does not give details on number of cross-validations)	Todorovic. (ECCV), 2012
73.70%	Leave One Group Out Cross-validation (25 cross-validations)	Solmaz, et al. (MVAP), 2012
72.60%	Leave One Group Out Cross-validation (25 cross-validations)	Kliper-Gross, et al. (ECCV), 2012

# Сравнение с другими моделями

## UCF50

UCF50 is an action recognition data set with 50 action categories, consisting of realistic videos taken from youtube. This data set is an extension of YouTube Action data set (UCF11) which has 11 action categories.

UCF50 data set's 50 action categories collected from youtube are: Baseball Pitch, Basketball Shooting, Bench Press, Biking, Biking, Billiards Shot, Breaststroke, Clean and Jerk, Diving, Drumming, Fencing, Golf Swing, Playing Guitar, High Jump, Horse Race, Horse Riding, Hula Hoop, Javelin Throw, Juggling Balls, Jump Rope, Jumping Jack, Kayaking, Lunges, Military Parade, Mixing Batter, Nun chucks, Playing Piano, Pizza Tossing, Pole Vault, Pommel Horse, Pull Ups, Punch, Push Ups, Rock Climbing Indoor, Rope Climbing, Rowing, Salsa Spins, Skate Boarding, Skiing, Skijet, Soccer Juggling, Swing, Playing Tabla, TaiChi, Tennis Swing, Trampoline Jumping, Playing Violin, Volleyball Spiking, Walking with a dog, and Yo Yo.

Source: [UCF50](#)

Homepage



License ⓘ

Unknown

Modalities

Videos

Languages

## Benchmarks

No benchmarks yet. [Start a new benchmark](#) or [link an existing one](#).

## Papers

Search for a paper or author

Paper	Code	Results	Date	Stars ↑
No matching records found				
Showing 0 to 0 of 0 entries				
		Previous	Next	

## Dataset Loaders

No data loaders found. You can [submit your data loader here](#).

## Tasks

## UCF101 (UCF101 Human Actions dataset)

Introduced by Soomro et al. in [UCF101: A Dataset of 101 Human Actions Classes From Videos in The Wild](#)

UCF101 dataset is an extension of UCF50 and consists of 13,320 video clips, which are classified into 101 categories. These 101 categories can be classified into 5 types (Body motion, Human-human interactions, Human-object interactions, Playing musical instruments and Sports). The total length of these video clips is over 27 hours. All the videos are collected from YouTube and have a fixed frame rate of 25 FPS with the resolution of 320 × 240.

Source: [Train-Driven Collaborative Learning with Spatial Temporal Attention for Video Classification](#)

Homepage



Source: <https://www.cse.cuhk.edu.hk/UCF/>

Usage



License ⓘ

Unknown

Modalities

Videos

Languages

## Benchmarks

Trend	Task	Dataset Variant	Best Model	Paper	Code
	Action Recognition	UCF101	SMART	<a href="#">Paper</a>	<a href="#">Code</a>
	Self-Supervised Action Recognition	UCF101	MEVideo	<a href="#">Paper</a>	<a href="#">Code</a>
	Zero-Shot Action Recognition	UCF101	MOV	<a href="#">Paper</a>	
	Video Frame Interpolation	UCF101	MA-CSPA	<a href="#">Paper</a>	
	Video Generation	UCF101	MCVD	<a href="#">Paper</a>	<a href="#">Code</a>

Show all 19 benchmarks

## Papers

Search for a paper or author

Paper	Code	Results	Date	Stars ↑
Spatiotemporal Contrastive Video Representation Learning	<a href="#">Code</a>	<a href="#">Results</a>	9 Aug 2020	74,845
AIVA: A Video Dataset of Spatio-temporally Localized Atomic Visual Actions	<a href="#">Code</a>	<a href="#">Results</a>	23 May 2017	74,843
Contextualized Spatio-Temporal Contrastive Learning with Self-Supervision	<a href="#">Code</a>	<a href="#">Results</a>	9 Dec 2021	74,843
VATT: Transformers for Multimodal Self-Supervised Learning from Raw Video, Audio and Text	<a href="#">Code</a>	<a href="#">Results</a>	22 Apr 2021	26,403
EVA: Exploring the Limits of Masked Visual Representation Learning at Scale	<a href="#">Code</a>	<a href="#">Results</a>	14 Nov 2022	22,459
Self-Supervised Multi-Modal Versatile Networks	<a href="#">Code</a>	<a href="#">Results</a>	29 Jun 2020	11,172
Learning Transferable Visual Models From Natural Language Supervision	<a href="#">Code</a>	<a href="#">Results</a>	26 Feb 2021	11,149
Depth-Aware Video Frame Interpolation	<a href="#">Code</a>	<a href="#">Results</a>	1 Apr 2019	7,652
Domain-Adversarial Training of Neural Networks	<a href="#">Code</a>	<a href="#">Results</a>	28 May 2015	5,512
A Large Scale Study on Unsupervised Spatiotemporal Representation Learning	<a href="#">Code</a>	<a href="#">Results</a>	29 Apr 2021	5,297

Showing 1 to 10 of 1,229 papers

Previous 1 2 3 4 5 ... 123 Next

## Dataset Loaders

<a href="#">pytorch/vision</a>	★ 12,938
<a href="#">pytorch/vision</a>	★ 12,938

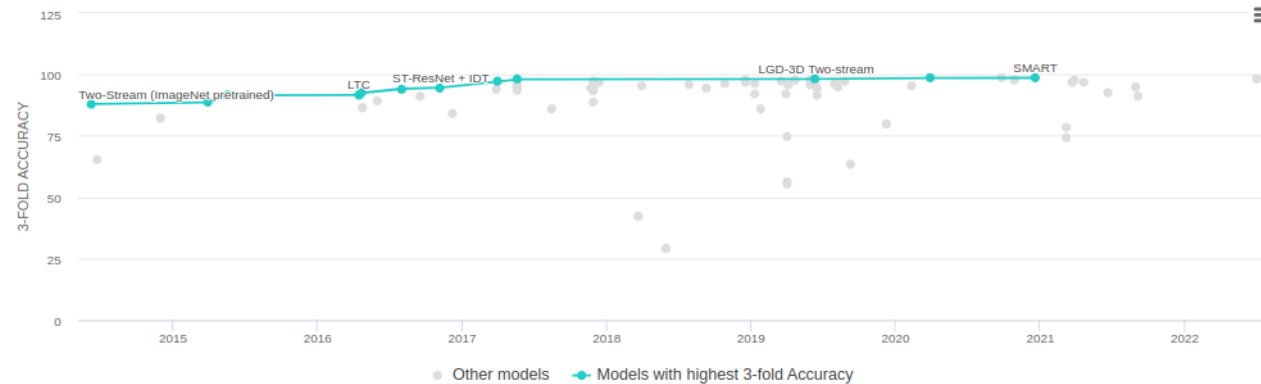


# Action Recognition on UCF101

Leaderboard

Dataset

View 3-fold Accuracy by Date for All models



Filter: LSTM ResNet GCN untagged

Edit Leaderboard

Rank	Model	3-fold Accuracy	Extra Training Data	Paper	Code	Result	Year	Tags
1	SMART	98.64	×	<a href="#">SMART Frame Selection for Action Recognition</a>			2020	
2	OmniSource (SlowOnly-8x8-R101-RGB + I3D-Flow)	98.6	✓	<a href="#">Omni-sourced Webly-supervised Learning for Video Recognition</a>			2020	
3	PERF-Net (multi-distilled S3D)	98.6	×	<a href="#">PERF-Net: Pose Empowered RGB-Flow Net</a>			2020	
4	LGD-3D Two-stream	98.2	×	<a href="#">Learning Spatio-Temporal Representation with Local and Global Diffusion</a>			2019	
5	Text4Vis	98.2	×	<a href="#">Revisiting Classifier: Transferring Vision-Language Models for Video Recognition</a>			2022	

# SMART на данный момент является SOTA для аналогичной задачи на UCF101

Method	Backbone	UCF101	HMDB51
Two-stream	VGG	92.5	62.4
I3D	Inc v3	98.0	80.7
DynaMotion + I3D	Inc v3	98.4	84.2
TSN	BN-Inc	94.2	69.9
KI-Net	Res-152	97.8	78.2
AAS	TSN	94.6	71.2
<b>SMART</b>	TSN	<b>95.8</b>	<b>74.6</b>
AAS	TSN+Kinetics	96.8	77.3
<b>SMART</b>	TSN+Kinetics	<b>98.6</b>	<b>84.3</b>

Table 5: Extending SMART to other approaches

Method	UCF101	HMDB51
ISTPAN	95.5	70.7
ISTPAN + SMART	<b>96.4</b>	<b>72.1</b>
I3D	98.0	80.0
I3D + Smart	<b>98.2</b>	<b>81.1</b>
STM-Resnet	94.2	68.9
STM-Resnet + SMART	<b>94.9</b>	<b>69.7</b>

# Проблемы

- Сравнительно маленький датасет UCF50, который входит в датасет побольше.
- Для работы с видео требуется значительно больше параметров
- Из-за отсутствия аугментаций модели не хватает информации, чтобы что-либо выучить