



How far are we from true AutoML

Winning solutions and results of *AutoDL challenge*

7th ICML AutoML Workshop on AutoML

Presented by Z. Liu in the name of
the AutoDL challenge team

The AutoDL challenge team

Original lead organizers:

- Olivier Bousquet (Google, Switzerland)
- André Elisseeff (Google, Switzerland)
- Isabelle Guyon (U. Paris-Saclay; UPSud/INRIA, France and ChaLearn, USA)
- Zhengying Liu (U. Paris-Saclay; UPSud, France)
- Wei-Wei Tu (4paradigm, China)

Other core team members:

- Sergio Escalera (U. Barcelona, Spain, and ChaLearn, USA)
- Julio Jacques Jr. (U. Barcelona, Spain)
- Meysam Madani (U. Barcelona, Spain)
- Adrien Pavao (U. Paris-Saclay; INRIA, France and ChaLearn, USA)
- Sebastien Treger (La Pallaise, France, and ChaLearn, USA)
- Zhen Xu (Ecole Polytechnique and U. Paris-Saclay; INRIA, France)

Other contributors to the organization, starting kit, and datasets, include:

- Stephane Ayache (AMU, France)
- Hubert Jacob Banville (INRIA, France)
- Mahsa Behzadi (Google, Switzerland)
- Kristin Bennett (RPI, New York, USA)
- Hugo Jair Escalante (IANOE, Mexico and ChaLearn, USA)
- Gavin Cawley (U. East Anglia, UK)
- Baiyu Chen (UC Berkeley, USA)
- Albert Clapes i Sintes (U. Barcelona, Spain)
- Bram van Ginneken (Radboud U. Nijmegen, The Netherlands)
- Alexandre Gramfort (U. Paris-Saclay; INRIA, France)
- Yi-Qi Hu (4paradigm, China)
- Tatiana Merkulova (Google, Switzerland)
- Shangeth Rajaa (BITS Pilani, India)
- Herilalaina Rakotoarison (U. Paris-Saclay, INRIA, France)
- Mehreen Saeed (FAST Nat. U. Lahore, Pakistan)
- Marc Schoenauer (U. Paris-Saclay, INRIA, France)
- Michele Sebag (U. Paris-Saclay; CNRS, France)
- Danny Silver (Acadia University, Canada)
- Lisheng Sun (U. Paris-Saclay; UPSud, France)
- Fengfu Li (4paradigm, China)
- Lichuan Xiang (4paradigm, China)
- Jun Wan (Chinese Academy of Sciences, China)
- Mengshuo Wang (4paradigm, China)
- Jingsong Wang (4paradigm, China)
- Ju Xu (4paradigm, China)

The challenge is running on the [Codalab platform](#), administered by [Université Paris-Saclay](#) and maintained by CKCollab LLC, with primary developers:

- Eric Carmichael (CKCollab, USA)
- Tyler Thomas (CKCollab, USA)

Sponsors:



Home institutions:



UNIVERSITAT DE
BARCELONA

Conferences:



ACML 2019

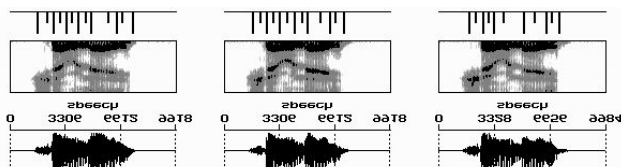
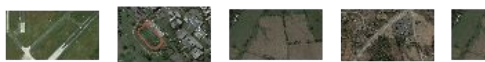
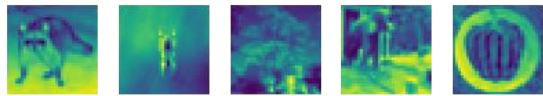


ECMLPKDD
Würzburg 2019

AutoDL Challenge Design

(1) Data: diverse modalities/domains

We formatted >100 datasets, 66 of which ended up being used in challenges



#	Dataset	Challenge(s)	Phase	Domain	Type	Class num.	Sample number	time	Tensor dimension	chal		
						train	test	row	col			
1	Mnist	AutoCV1	public	HWR	image	10	60000	10000	1	28	28	1
2	Cmn5v	AutoCV1	public	objects	image	100	40661	11819	1	32	32	3
3	Pedro	AutoCV1	public	people	image	20	80095	19905	1	var	var	3
4	Usps10	AutoCV1	public	aerial	image	11	634	166	1	var	var	3
5	Mnist	AutoCV1	public	medical	image	7	8000	1900	1	var	var	3
6	Ukaidle	AutoCV1	feedback	HWR	image	3	6979	1719	1	var	var	3
7	Kwiazek	AutoCV1	feedback	medical	image	20	23118	6269	1	var	var	3
8	Beatsix	AutoCV1	feedback	people	image	15	4406	1094	1	350	350	3
9	Serius	AutoCV1	final	aerial	image	3	32400	8160	1	28	28	1
10	Hippocrate	AutoCV1	feedback	medical	image	2	17917	41108	1	96	96	3
11	Lonsome	AutoCV2	final	HWR	image	3	27938	6939	1	var	var	3
12	Tim	AutoCV1	final	objects	image	200	80000	20000	1	32	32	3
13	Spelion	AutoCV2	final	people	image	100	6077	1514	1	var	var	3
14	Ideal	AutoCV1	final	aerial	image	45	25231	6269	1	256	256	3
15	Ray	AutoCV1	final	medical	image	7	4492	1114	1	976	976	3
16	Kraut	AutoCV2	public	action	video	4	1528	863	var	120	160	1
17	Kraut	AutoCV2	public	action	video	6	1528	863	var	120	160	1
18	Kraut	AutoCV2	public	action	video	4	1528	863	var	60	80	1
19	Freddy	AutoCV2	feedback	HWR	image	2	54056	13631	1	var	var	3
20	Hemur	AutoCV2	feedback	action	video	12	1354	333	var	var	var	3
21	Inac2	AutoCV2	feedback	action	video	249	38372	9501	var	102	78	1
22	Fernax	AutoCV2	final	action	video	4	32994	8203	var	80	80	3
23	Fiona	AutoCV2	final	action	video	6	8038	1562	var	var	var	3
24	Monac1	AutoCV2	final	action	video	20	10380	2565	var	168	168	3
25	Kitouse	AutoCV2	final	action	video	25	18602	4963	var	66	82	3
26	data01	AutoSpeech	public	speech	time	100	3000	3000	var	1	1	1
27	data02	AutoSpeech	public	speech	time	7	428	107	var	1	1	1
28	data03	AutoSpeech	public	speech	time	4	48	290	var	1	1	1
29	data04	AutoSpeech	public	speech	time	20	579	274	var	1	1	1
30	data05	AutoSpeech	public	speech	time	10	199	507	var	1	1	1
31	data11	AutoSpeech	feedback	speech	time	50	1300	2000	var	1	1	1
32	data12	AutoSpeech	feedback	speech	time	5	3120	346	var	1	1	1
33	data13	AutoSpeech	feedback	speech	time	3	2000	1330	var	1	1	1
34	data14	AutoSpeech	feedback	speech	time	8	767	191	var	1	1	1
35	data15	AutoSpeech	feedback	speech	time	76	2566	791	var	1	1	1
36	data21	AutoSpeech	final	speech	time	50	800	1200	var	1	1	1
37	data22	AutoSpeech	final	speech	time	8	49	291	var	1	1	1
38	data23	AutoSpeech	final	speech	time	3	2000	284	var	1	1	1
39	data24	AutoSpeech	final	speech	time	16	384	386	var	1	1	1
40	data25	AutoSpeech	final	speech	time	100	3008	752	var	1	1	1
41	O1	AutoNLP	public	english	text	2	7792	1821	var	1	1	1
42	O2	AutoNLP	public	english	text	20	13174	7522	var	1	1	1
43	O3	AutoNLP	public	english	text	2	60000	40000	var	1	1	1
44	O4	AutoNLP	public	chinese	text	2	50000	10000	var	1	1	1
45	O5	AutoNLP	public	chinese	text	18	150000	22000	var	1	1	1
46	PU1	AutoNLP	feedback	english	text	492	492	492	var	1	1	1
47	PU2	AutoNLP	feedback	english	text	5	13851	23499	var	1	1	1
48	PU3	AutoNLP	feedback	chinese	text	2	110203	19519	var	1	1	1
49	PU4	AutoNLP	feedback	chinese	text	11	100000	50000	var	1	1	1
50	PU5	AutoNLP	feedback	chinese	text	31	600000	400000	var	1	1	1
51	PR1	AutoNLP	final	english	text	20	43007	6907	var	1	1	1
52	PR2	AutoNLP	final	english	text	2	42500	7501	var	1	1	1
53	PR3	AutoNLP	final	english	text	4	90000	30050	var	1	1	1
54	PR4	AutoNLP	final	chinese	text	11	130000	50000	var	1	1	1
55	PR5	AutoNLP	final	chinese	text	15	250000	132888	var	1	1	1
56	AutoL	AutoDL	public	categorical	tabular	5	39074	6768	1	1	24	1
57	Dilbert	AutoDL	public	objects	tabular	5	14860	9720	1	1	2000	1
58	Image	AutoDL	public	HWR	tabular	1	32600	32600	1	1	1569	1
59	Maxline	AutoDL	public	artificial	tabular	2	4220	3240	1	1	259	1
60	Beats	AutoDL	CE pair	audio	tabular	4	2160	2160	1	1	1	1
61	Bliss	AutoDL	final	audio	tabular	20	10931	2733	1	1	400	1

- IMAGE
- VIDEO
- SPEECH
- TEXT
- TABULAR
- Multi-label tasks

Liu Z, Xu Z, Rajaa S, Madadi M. Towards Automated Deep Learning: Analysis of the AutoDL challenge series 2019. To appear in *NeurIPS CD2019* in Proceedings of Machine Learning Research (PMLR) 2019:10.

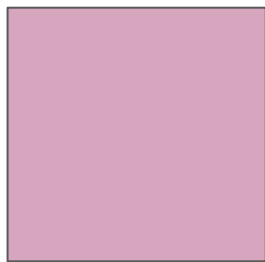
(1) Data: RAW data in a Tensor Format

dim (1, 1, y, 1)

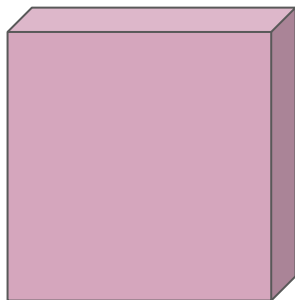
Feature vectors
(tabular data)



dim (1, x, y, 1)



dim (1, x, y, c)



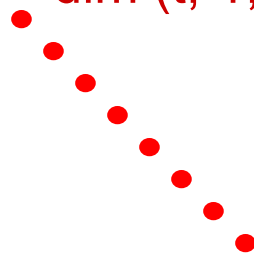
Images

dim (t, 1, 1, 1)



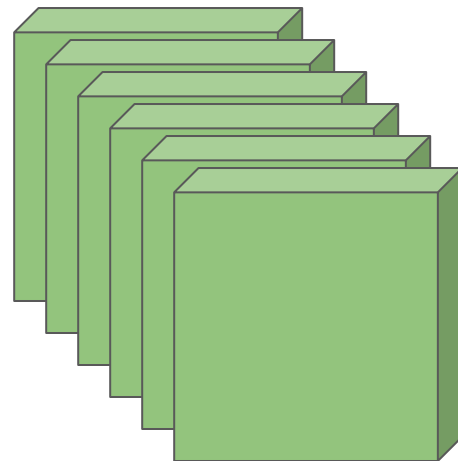
Speech

dim (t, 1, 1, 1)



Text (with a vocabulary as meta-data)

dim (t, x, y, c)



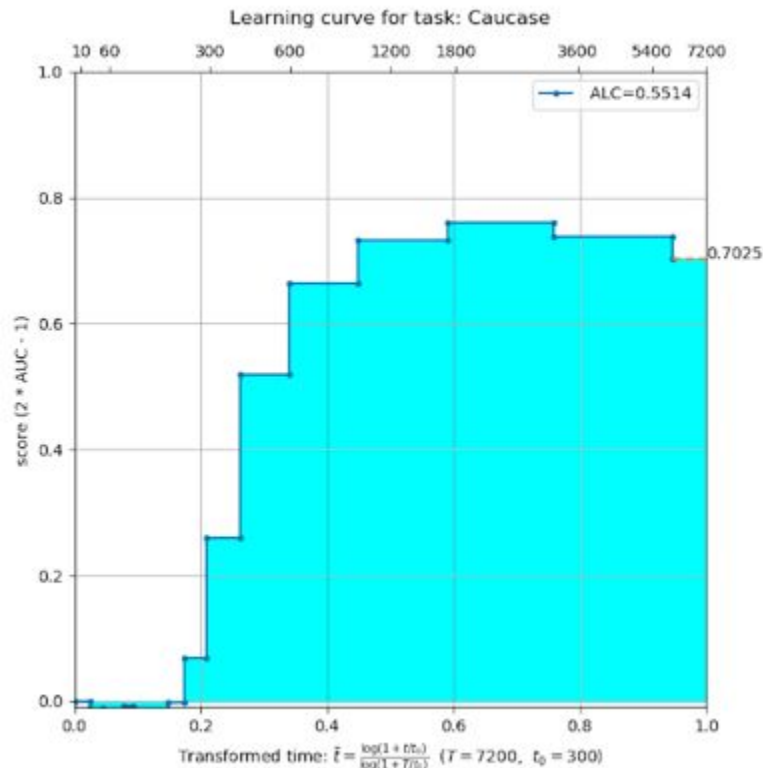
Videos

(2) Evaluation: Fixed max T + Any-time Learning

Time rescaling:

$$\tilde{t}(t) = \frac{\log(1 + t/t_0)}{\log(1 + T/t_0)}$$

$$\begin{aligned} ALC &= \int_0^1 s(t) d\tilde{t}(t) \\ &= \int_0^T s(t) \tilde{t}'(t) dt \\ &= \frac{1}{\log(1 + T/t_0)} \int_0^T \frac{s(t)}{t + t_0} dt \end{aligned}$$



(2) Evaluation: Blind testing

Two phases:

- **Feed-back phase:** 5 datasets, 5 submissions/day for 3-4 months.
- **Final test phase:** 10 OTHER unseen datasets, **ONE** single submission.

**BOTH phases, code TRAINED and TESTED
on the platform => datasets invisible.**

Additional **“public” datasets**
=> **META-LEARNING.**



(3) Starting Kit and Baselines

Baseline 0: **Constant** predictions (for debug purposes)

domain-agnostic

Baseline 1: **Linear model**

Baseline 2: **Multi-dimensional CNN** (auto-rescaling)

domain-**dependent**

Baseline 3: Combination of winning solutions from previous challenges:

Image & Video: Kakaobrain, ResNet (He et al, 2016) and Fast Auto Augment (Cubuk et al. (2018); Lim et al. (2019))

Speech: PASA_NJU, Spectral transform, logistic reg., lightGBM, CNN, ResNet, VggVox, LSTM, etc.

Text: Upwind_flys, LinearSVC, LSTM, BERT, etc., selected with meta-controller

Tabular (new): Fully connected network.

Challenge Design Recap

- (1) Raw data from 5 domains: Image, Video, Speech, Text, Tabular.
- (2) Fixed time budget. Any-time learning (ALC metric). Blind testing.
- (3) Starting kit, sample “public” data and baselines provided.
- (4) Fixed computational resources.
- (5) Using Deep Learning was NOT imposed.

Challenge Results



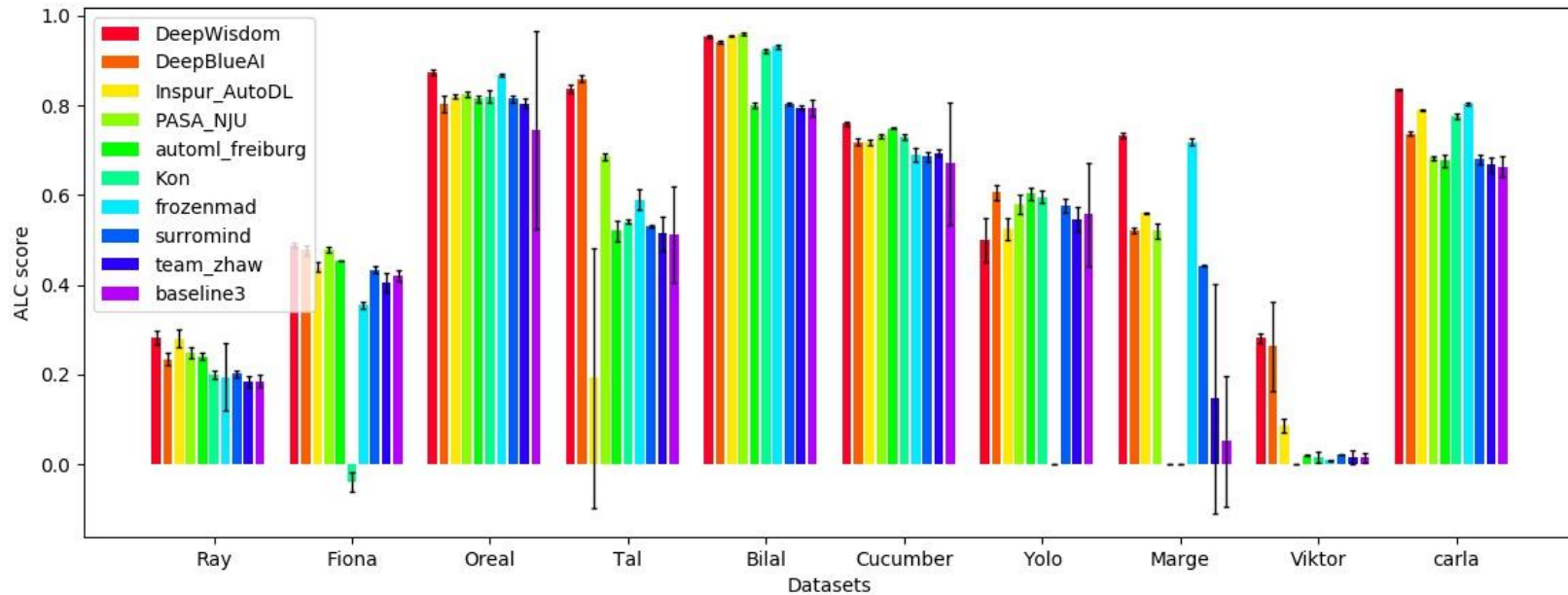
WINNERS



(code open-sourced at autodl.chalearn.org)

Challenge	1st Prize (\$2000)	2nd Prize (\$1500)	3rd Prize (\$500)
AutoCV	kakaobrain (Kakao Brain)	DKKimHCLee (Hana. Tech. Inst.)	base_1 (Hanyang University)
AutoCV2	kakaobrain (Kakao Brain)	tanglang (Xiamen University)	kvr (-)
AutoNLP	DeepBlueAI (DeepBlue Technology)	upwind_fly (Lenovo)	txta (gsdata.cn)
AutoSpeech	PASA_NJU (Nanjing University)	DeepWisdom (fuzhi.ai)	Kon (NS Solutions Corporation)
AutoDL	DeepWisdom (fuzhi.ai)	DeepBlueAI (DeepBlue Technology)	Inspur_AutoDL & PASA_NJU

AutoDL final phase results



Did we get answers to our questions?

(1) Unified approach for ALL 5 domains?

(Image, Video, Speech, Text, Tabular.)

(2) Time budget sufficient? Any-time learning possible?

(3) Was sample “public” data used for meta-learning?

(1) Unified approach for ALL 5 domains?

Participant survey:

[DOMAIN-DEPENDENCY] Do you use a separate approach for each dataset/domain (rather than a single unified approach)

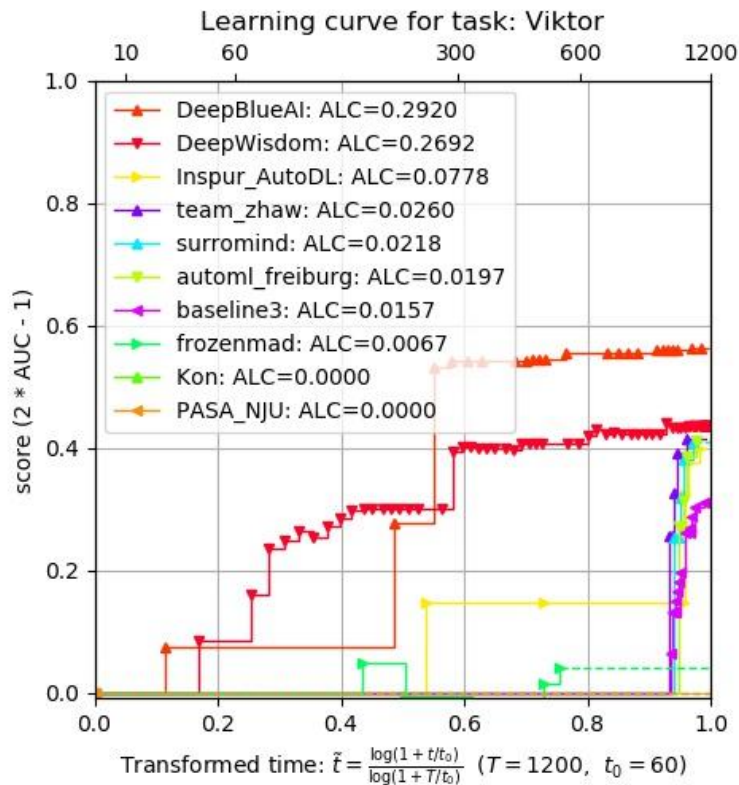
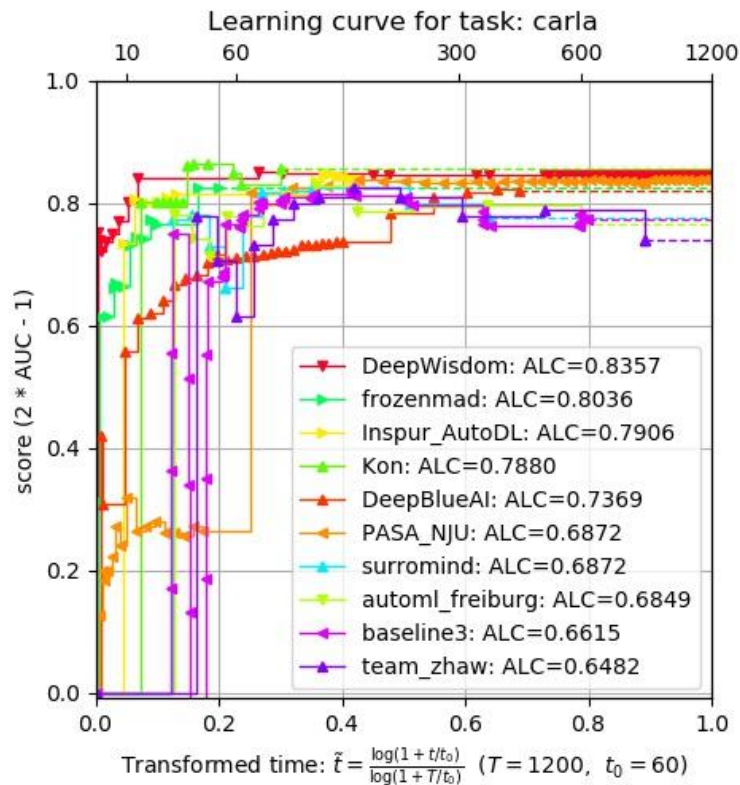
10 responses

100%



● Yes
● No

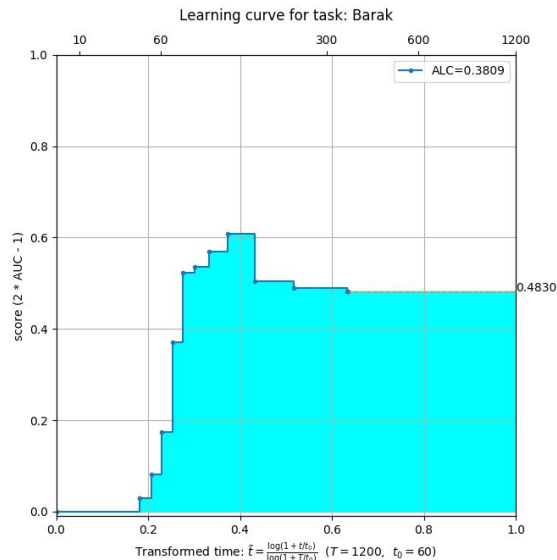
(2) Time budget sufficient? Any-time learning possible?



Impact of t_0

$$\tilde{t}(t) = \frac{\log(1 + t/t_0)}{\log(1 + T/t_0)}$$

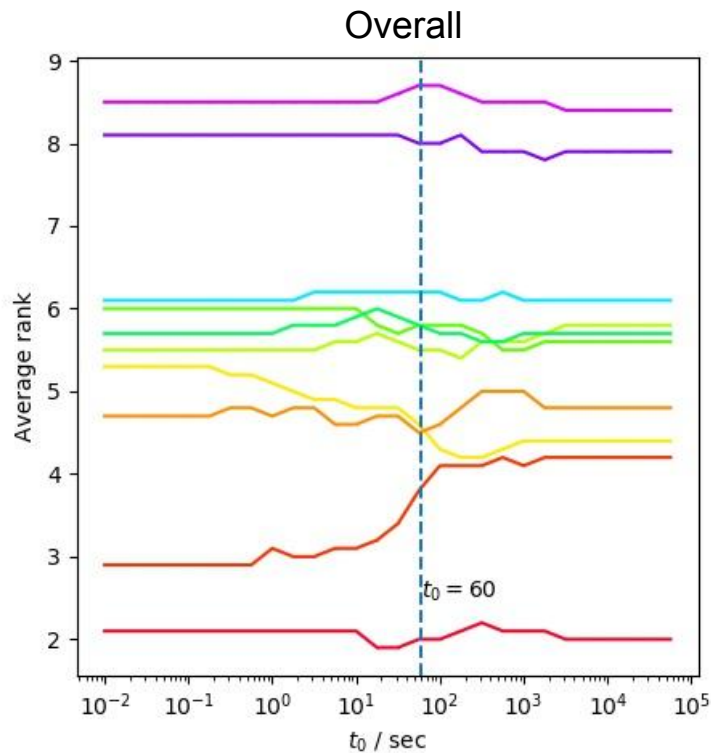
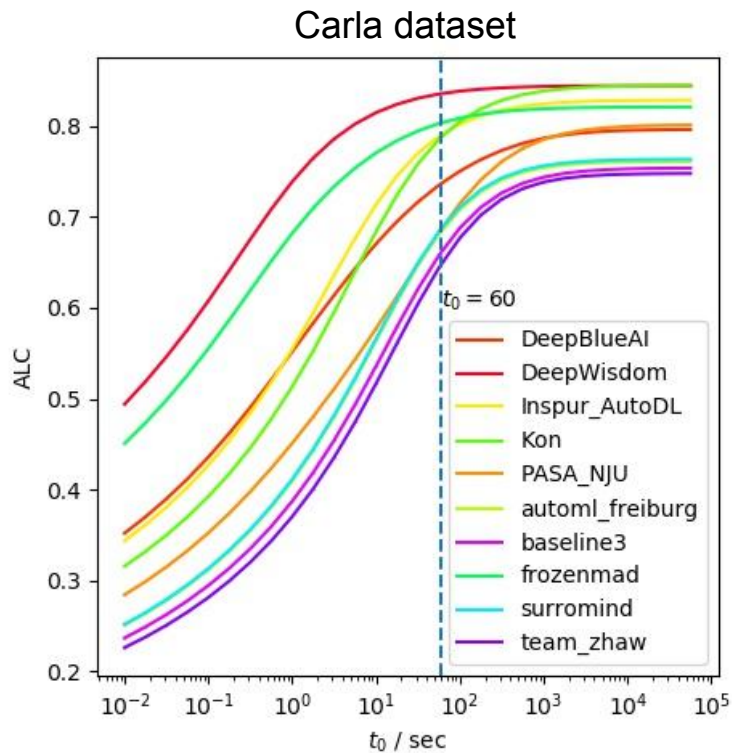
$$\begin{aligned} ALC &= \int_0^1 s(t) d\tilde{t}(t) \\ &= \int_0^T s(t) \tilde{t}'(t) dt \\ &= \frac{1}{\log(1 + T/t_0)} \int_0^T \frac{s(t)}{t + t_0} dt \end{aligned}$$



$$\lim_{t_0 \rightarrow 0^+} ALC(t_0) = s(0)$$

$$\lim_{t_0 \rightarrow +\infty} ALC(t_0) = \frac{1}{T} \int_0^T s(t) dt$$

Impact of t_0

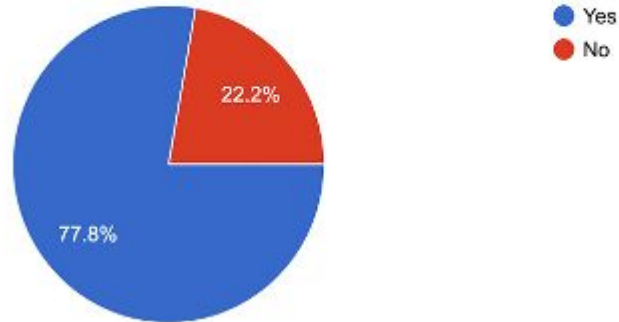


(3) Was “public” data used for meta-learning?

Participant survey:

[META-LEARNING] Did you use the public datasets (or other data available to you) for model selection or apply meta-learning techniques?

9 responses

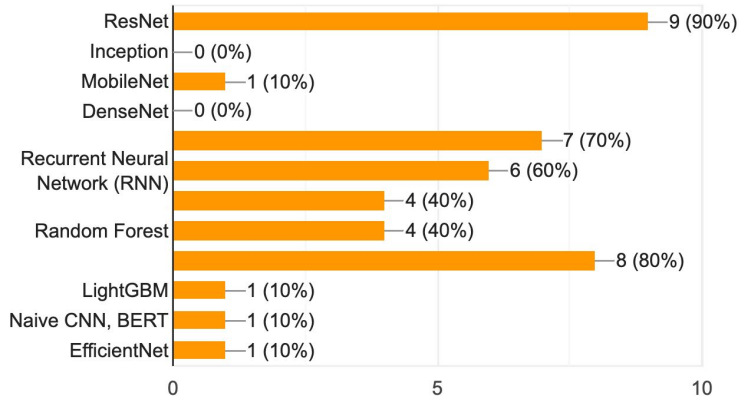


Winning solutions

Neural architectures used in the winning approaches

Base predictor / architecture

10 responses



Architecture name	# Parameters	Domains	Teams
ResNet-18, ResNet-9 (He et al 2015)	11.4M, 5.7M	image, video	Kakaobrain, DeepWisdom, automl_freiburg
MC3 (Du Tran et al CVPR 2018)	32.8M	video	DeepWisdom
EfficientNet-(b0, b1, b2) (M. Tan and Q. Le. 2019)	5.3M, 7.8M, 9.2M	image, video	DeepWisdom, automl_freiburg
MobileNetV2 (M. Sandler et al 2019)	3.4M	image, video	team_zhaw, DeepBlueAI
TextCNN	variable	text	Upwind_flys, DeepWisdom
Fast RCNN (Ross Girshick)		text	DeepWisdom
LSTM, BiLSTM (Hochreiter, Schmidhuber 1997)	0.2M-1M	text, speech	frozenmad, PASA_NJU
GRU, BiGRU, (Kyunghyun Cho et al 2014) GRU with Attention	0.1M-1M	text, speech	DeepBlueAI, DeepWisdom
BERT-like (Tiny-BERT(X.Jiao et al))	<110M	text	frozenmad, upwind_flys
DNN	<1M	tabular	DeepWisdom

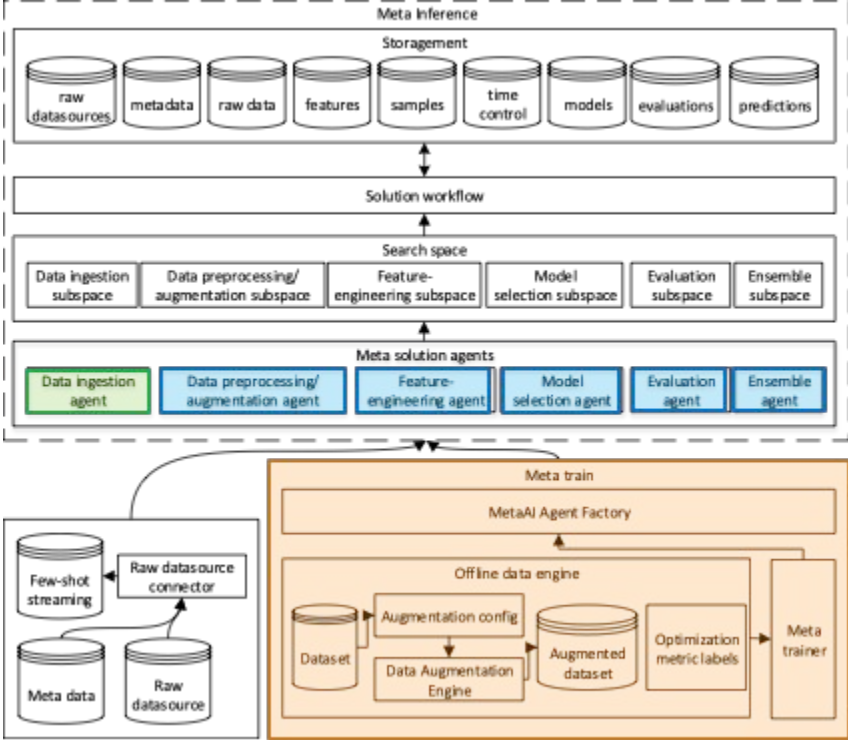
AutoML techniques vs domains

Approach	image	video	speech	text	tabular
Meta-learning	Offline meta-training transferred with AutoFolio [25] based on meta-features (<i>automl freiburg</i>) Offline meta-training generating solution agents, searching for optimal sub-operators in predefined sub-spaces, based on dataset meta-data. (<i>DeepWisdom</i>) MAML-like method [17] (<i>team zhaw</i>)				
Preprocessing	image cropping and data augmentation (<i>PASANJU</i>), fast autoaugment (<i>DeepBlueAI</i>)	Sub-sampling keeping 1/6 frames and adaptive image size (<i>DeepBlueAI</i>) Adaptive image size	MFCC, Mel Spectrogram, STFT	root features extractions with stemmer, meaningless words filtering (<i>DeepBlueAI</i>)	Numerical and Categorical data detection and encoding
Hyperparameter Optimization	Offline with BOHB [26] (Bayesian Optimization and Multi-armed Bandit) (<i>automl freiburg</i>) Sequential Model-Based Optimization for General Algorithm Configuration (SMAC) (<i>automl freiburg</i>)				Bayesian Optimization (<i>PASANJU</i>) HyperOpt [27] (<i>Inspur AutoDL</i>)
Transfer learning	Pre-trained on ImageNet [28] (all teams except <i>Kon</i>)	Pre-trained on ImageNet [28] (all top-8 teams except <i>Kon</i>) MC3 model pretrained on Kinetics (<i>DeepWisdom</i>)	ThinResnet34 pre-trained on VoxCeleb2 (<i>DeepWisdom</i>)	FastText pre-trained on Common Crawl (<i>frozenmad</i>)	
Ensemble learning	Adaptive Ensemble Learning (ensemble latest 2 to 5 predictions) (<i>DeepBlueAI</i>)	Ensemble Selection [29] (top 5 validation predictions are fused) (<i>DeepBlueAI</i>); Ensemble models sampling 3, 10, 12 frames (<i>DeepBlueA</i>)	last best predictions ensemble strategy (<i>DeepWisdom</i>) averaging 5 best overall and best of each model: LR, CNN, CNN+GRU (<i>DeepBlueA</i>)	Weighted Ensemble over 20 best models [29] (<i>DeepWisdom</i>)	LightGBM ensemble with bagging method [30] (<i>DeepBlueAI</i>), Stacking and blending (<i>DeepWisdom</i>)

Teams vs domains

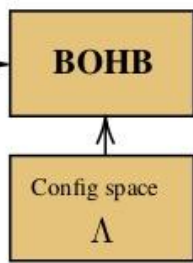
team	image	video	speech	text	tabular
DeepWisdom	[ResNet-18 and ResNet-9 models] [pretrained on ImageNet]	[MC3 model] [pretrained on Kinetics]	[fewshot learning] [LR, ThinResnet34 models] [pretrained on VoxCeleb2]	[fewshot learning] [task difficulty and similarity evaluation for model selection] [SVM, TextCNN , [fewshot learning] RCNN, GRU, GRU with Attention]	[LightGBM , Xgboost, Catboost, DNN models] [no pretrained]
DeepBlueAI	[data augmentation with Fast AutoAugment] [ResNet-18 model]	[subsampling keeping 1/6 frames] [Fusion of 2 best models]	[iterative data loader (7, 28, 66, 90%)] [MFCC and Mel Spectrogram preprocessing] [LR, CNN, CNN+GRU models]	[Samples truncation and meaningless words filtering] [Fasttext, TextCNN , BiGRU models] [Ensemble with restrictive linear model]	[3 LightGBM models] [Ensemble with Bagging]
PASA NJU	ResNet-18 and SeResnext50; preprocessing: shape standardization and image flip (data augmentation)	ResNet-18 and SeResnext50; preprocessing: shape standardization and image flip (data augmentation)	[data truncation(2.5s to 22.5s)] [LSTM, VggVox ResNet with pretrained weights of DeepWis- dom(AutoSpeech2019) ThinResnet34?]	[data truncation(300 to 1600 words)] [TF-IDF and word embedding]	[iterative data loading] [Non Neural Nets models] [models complexity increasing over time] [Baysien Optimization of hyperparameters]
frozenmad	[images resized under 128x128] [progressive data loading increasing over time and epochs] [ResNet-18 model] [pretrained on ImageNet]	[Successive frames difference as input of the model] [pretrained ResNet-18 with RNN models]	[progressive data loading in 3 steps 0.01, 0.4, 0.7] [time length adjustment with repeating and clipping] [STFT and Mel Spectrogram preprocessing] [LR, LightGBM, VggVox models]	[TF-IDF and BERT tokenizers] [SVM, RandomForest, CNN, tinyBERT]	[progressive data loading] [no preprocessing] [Vanilla Decision Tree, RandomForest, Gradient Boosting models applied sequentially over time]

DeepWisdom

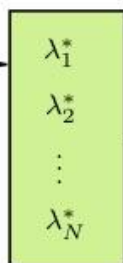


automl_freiburg

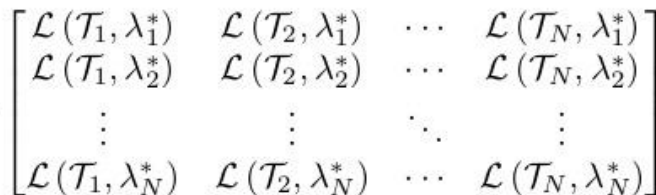
Meta-train tasks



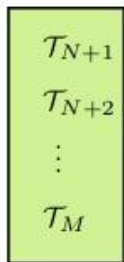
Optimal hyperparameters



Performance matrix



Meta-test tasks

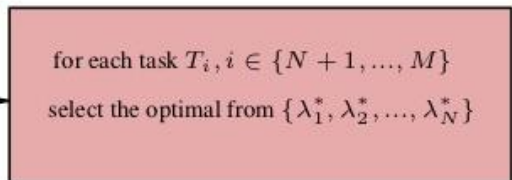


Meta-features

Meta-features



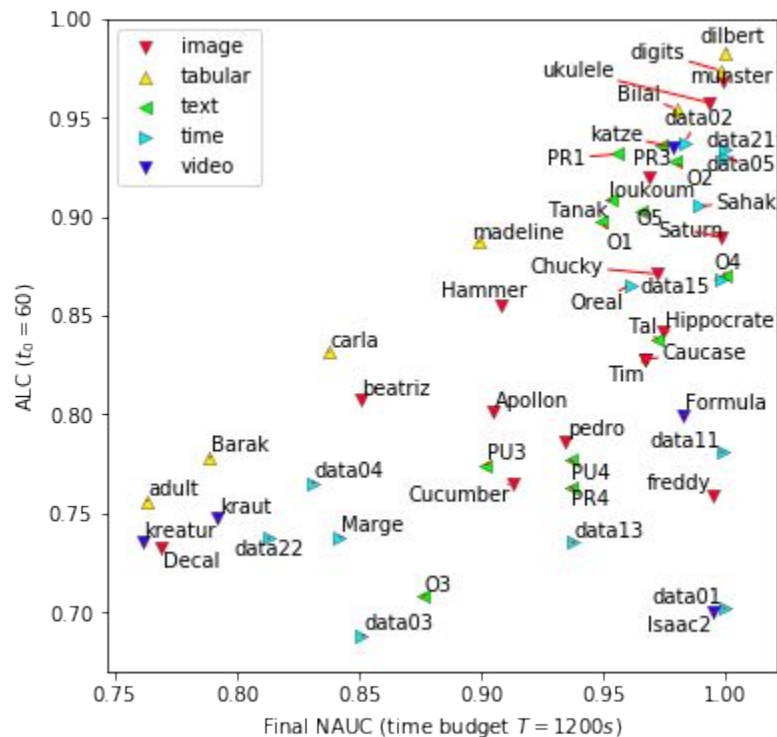
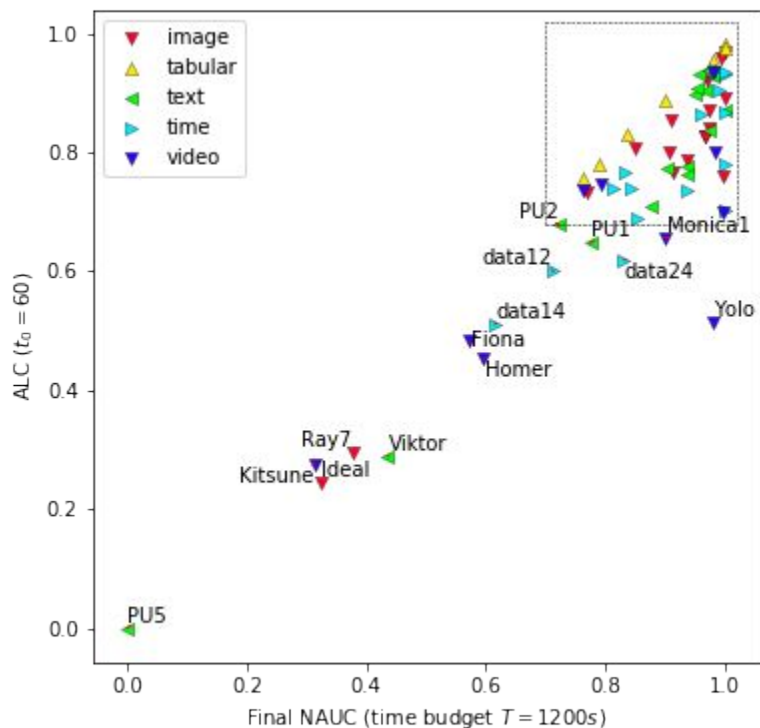
Select



Winning Solutions Recap

- (1) **Deep learning** is still dominant;
- (2) **Fixed domain-dependent** and/or **pre-trained** neural architectures are heavily used
- (3) **Neural architecture search** (NAS) hasn't been employed due to its huge computational cost
- (4) **Meta-learning** and **data loading/ingestion** strategies are used (and are useful)

Benchmarking: DeepWisdom on all 66 AutoDL datasets



More info at: autodl.chalearn.org

References

- [1] Liu Z, Bousquet O, Elisseeff A, et al. AutoDL Challenge Design and Beta Tests-Towards automatic deep learning. In: *MetaLearn Workshop @ NeurIPS2018*. Montreal, Canada; 2018. <https://hal.archives-ouvertes.fr/hal-01906226>. Accessed October 2, 2019.
- [2] Liu Z, Guyon I, Junior JJ, et al. AutoCV Challenge Design and Baseline Results. In: *CAP 2019 - Conférence Sur l'Apprentissage Automatique*. Toulouse, France; 2019. <https://hal.archives-ouvertes.fr/hal-02265053>. Accessed November 5, 2019.
- [3] Liu Z, Xu Z, Escalera S, et al. Towards Automated Computer Vision: Analysis of the AutoCV Challenges 2019. To appear in *Pattern Recognition Letters* of Elsevier. 2020. <https://hal.archives-ouvertes.fr/hal-02386805>. Accessed December 6, 2019.
- [4] Liu Z, Xu Z, Rajaa S, Madadi M. Towards Automated Deep Learning: Analysis of the AutoDL challenge series 2019. To appear in *NeurIPSCD2019* in Proceedings of Machine Learning Research (PMLR) 2019:10.
- [5] Liu Z, et al, Post-challenge analysis of AutoDL challenges 2019, submitted to TPAMI.

Takeaways

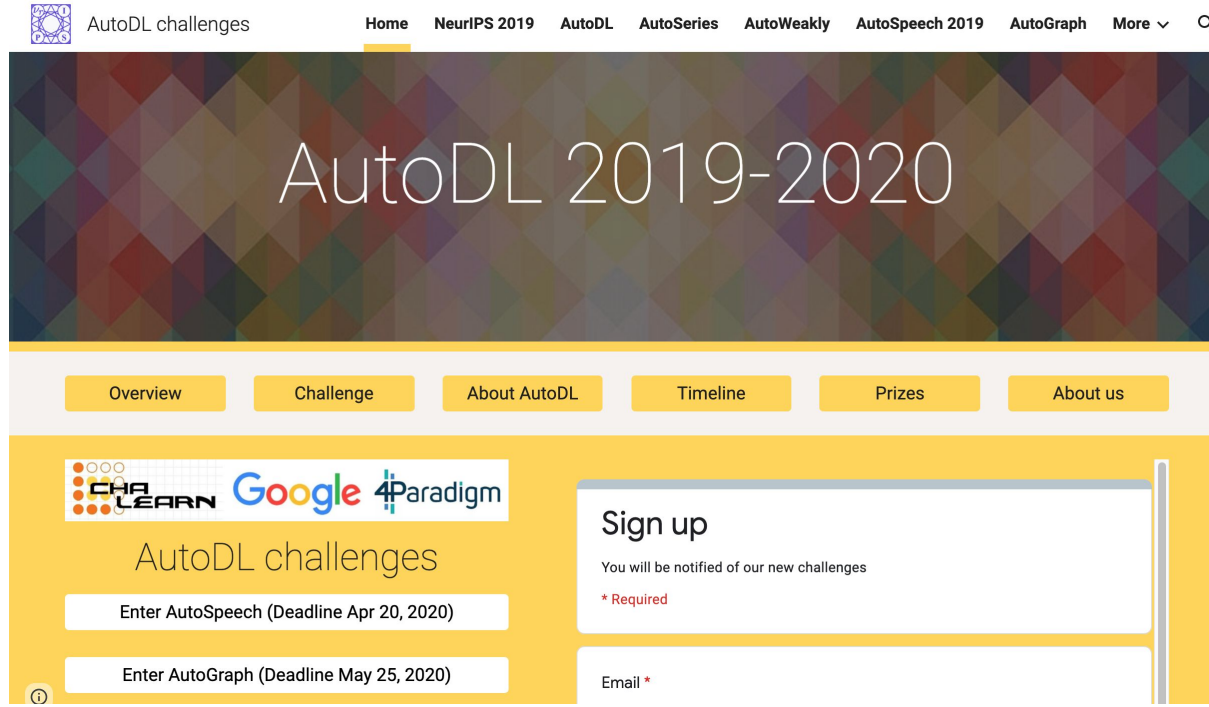
Lessons learned

- (1) The winning methods are capable of generalizing on new unseen datasets => **Potential universal AutoML solutions**
- (2) Domain-dependent approaches are dominant
=> **No universal workflows, mostly hand-tuned meta-learning**
- (3) We cannot afford to run expensive NAS for every new task
=> **Need transferability of learned architectures**
- (4) Beating Baseline 3 by using “true” meta-learning is hard
=> **Need more meta-train datasets (public datasets)**

To achieve true AutoML, we need...

- (1) Constructive and efficient **representation of meta-knowledges**: domain/modality related, pixel correlation, etc
- (2) Constructive and efficient **representation of learning algorithms**: architecture encoding, code-based, etc
- (3) **Transferable** neural architecture search (NAS) to learn a fast algorithm/function: meta-knowledges -> architecture
- (4) **Lifelong learning** systems and/or **world models** that can learn ONCE but continuously

Thank you! Questions?



The screenshot shows the homepage of the AutoDL challenges website. At the top left is the logo for AutoDL challenges, a square with a grid of numbers. To its right is the text "AutoDL challenges". The navigation menu includes "Home", "NeurIPS 2019", "AutoDL", "AutoSeries", "AutoWeekly", "AutoSpeech 2019", "AutoGraph", and "More" with a dropdown arrow and a search icon. The main banner features a colorful geometric pattern and the text "AutoDL 2019-2020". Below the banner is a row of yellow buttons: "Overview", "Challenge", "About AutoDL", "Timeline", "Prizes", and "About us". The main content area has a yellow background. On the left, there are logos for CHA LEARN, Google, and 4Paradigm. Below them is the text "AutoDL challenges" and two input fields: "Enter AutoSpeech (Deadline Apr 20, 2020)" and "Enter AutoGraph (Deadline May 25, 2020)". On the right, there is a "Sign up" section with the text "You will be notified of our new challenges" and a red asterisk indicating a required field. Below this is an "Email" input field with a red asterisk.

JOIN THE CHALLENGES!

autodl.chalearn.org