

# A. Appendix for YieldSAT: A Multimodal Benchmark Dataset for High-Resolution Crop Yield Prediction

## A.1. Crop Values for Data Preprocessing

The maximum accepted yield values for each crop type are given in Tab. 6. Yield values above this threshold were removed. Additionally, the standard moisture is given in Tab. 6 for each crop type that is used to calculate the scaled yield as shown in the main paper.

## A.2. Data Analysis

### A.2.1. Data Quality and Support Size

The rasterization workflow for a single field is shown in Fig. 11. Combine harvester in vector format must be processed to match high-resolution EO inputs (e.g., S2). For this, every yield measurement that falls into a grid cell of  $10\text{ m} \times 10\text{ m}$  are averaged. Due to harvester path density, swath width, speed, logging frequency, and positional delay, different numbers of yield measurements can fall into a single grid cell, resulting in variable support sizes for every pixel. Those support sizes are often correlated in space. An example of different support sizes for every yield map quality level, shown in Fig. 5, is illustrated in Fig. 12. The figure shows the number of points that fall into a grid cell and the corresponding standard deviation per pixel. It is evident that the yield maps have different support sizes and variability per pixel. In Fig. 13, we show that these characteristics depend on the crop type and dataset. For instance, Brazil

Table 6. Maximum accepted yield values in t/ha for every crop type and the standard moisture.

Standard Values	Wheat	Rapeseed	Soybean	Corn
Max. Yield (t/ha)	20	10	15	45
Standard Moisture (%)	15	9	15	16

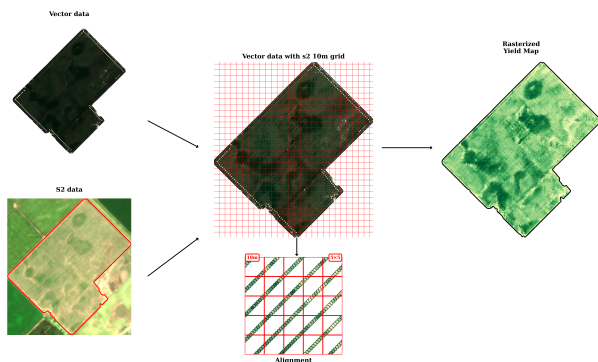


Figure 11. Schematic overview of data collection and preprocessing: A combine harvester collects point vector data of the yield, which is preprocessed and rasterized to align with S2 10 m grid for pixel-wise regression.

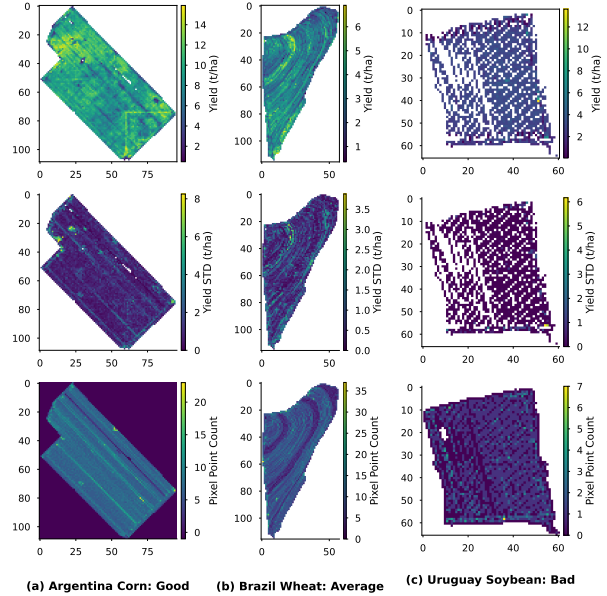


Figure 12. Label support size and standard deviation per 10 m pixel caused by the rasterization. Top row: sample count, bottom row: standard deviation.

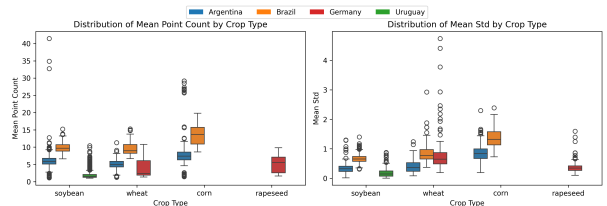


Figure 13. Boxplots for pixel point count (left) and mean standard deviation (right) per crop and country.

has consistently more yield points per pixel, while Uruguay has fewer samples. This is also reflected in a lower standard deviation per pixel. For each yield map, we also provide the raster grid for the variable support size.

Each field was manually inspected and curated. For this, each field is assigned a quality score of either "good", "average", or "bad", following a stringent guideline:

**1. Good** yield maps are characterized by high data quality and consistency. They do not exhibit strong visible striping artifacts (small can be present). The data coverage is dense, with no significant gaps or sparsity. Additionally, good yield maps show meaningful variability across the field, reflecting realistic yield differences, while maintaining smooth spatial transitions without abrupt or isolated changes.

**2. Bad** yield maps display clear quality issues. These may include visible striping patterns, often caused by sensor or machinery inconsistencies, as well as sparse or incomplete data coverage. Such maps typically lack meaningful vari-

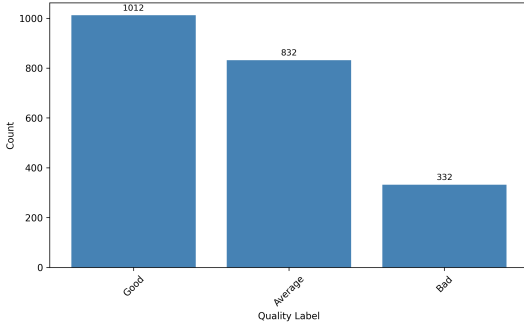


Figure 14. Overall distribution of yield map quality over the entire dataset.

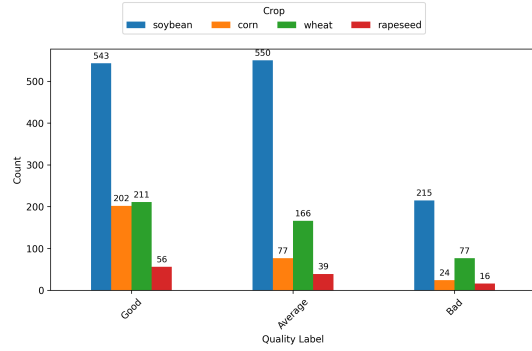


Figure 16. Overall distribution of yield map quality over the entire dataset grouped by country.

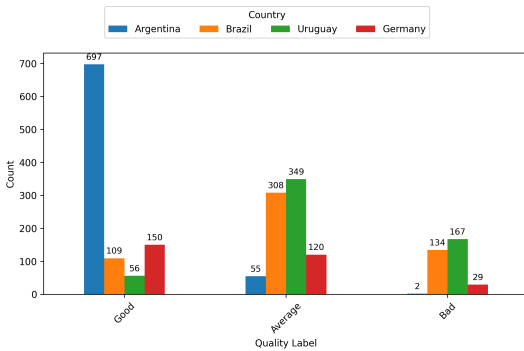


Figure 15. Overall distribution of yield map quality over the entire dataset grouped by country.

ability, with large areas showing nearly constant values. They may also contain isolated pixels or small regions with sudden, unrealistic changes in yield values that are inconsistent with the surrounding data.

**3. Average** yield maps fall between these two extremes. They may exhibit minor striping, moderate sparsity, or limited variability, but not to the extent observed in bad maps. While generally usable, they do not meet the quality standards of good yield maps and may require additional processing or careful interpretation.

An example of each class is shown in the main paper in Fig. 5. The distribution of yield map quality across the entire dataset is shown in Fig. 14. The plot indicates that the majority of yield maps are of good quality, followed by average and bad categories, based on a total of 332 fields. Fig. 15 presents the distribution of yield map quality grouped by country. Argentina exhibits the highest number of good-quality yield maps, whereas Uruguay has the largest number of low-quality (bad) yield maps. Finally, Fig. 16 shows the distribution of yield map quality by crop type. Soybean accounts for the largest number of high-quality yield maps, followed by wheat and corn.

The quality labels are stored in the metadata of each field.

### A.2.2. Distribution Shifts: Countries and Crop Types

To test whether the distributions of crop types and countries are significantly different, significance testing is performed to compare the distributions between groups. Both tests demonstrate significantly different distributions between countries and crop types ( $p$ -values  $< 0.0001$ ) as shown in Tab. 7. Following, we perform a post-hoc analysis that compares pairs within each group using Dunn’s significance test using a Holm-Bonferroni correction. The results for the pairwise comparison between each country are shown in Tab. 8. Notably, each group has a distinct yield distribution, with most  $p$ -values  $< 0.0001$ . Except for between Germany and Brazil, no significance is present. Additionally, the Dunn’s test for the pairwise comparison between crop types is shown in Tab. 9. Similarly, most distributions differ significantly between crop types, with all  $p$ -values  $< 0.0001$ , except for rapeseed and soybean.

### A.2.3. Distribution Shifts: Regions and Years

Importantly, we also find significant differences in data distribution and surface reflectance across regions and years within a single country and crop type. In Fig. 17 t-SNE plot of the surface reflectance for only soybean in Argentina is shown. The data are colored by year (left) and region (right). We highlight that, in both cases, a separation between individual years and regions is evident, as shown by clusters within each group. Below is the kernel density estimation plot of the target yield distribution. The data are grouped by year (left) and region (right). Note that each group’s distribution is unique. Individual distributions do not follow a normal distribution. Moreover, the plot shows that the yield distribution differs between years and regions. The results indicate unique patterns within each year-region combination across countries and crop types.

To test the hypothesis that all groups have the same yield distribution, a Kruskal-Wallis test is performed to compare the distributions across groups. The results are depicted in Tab. 10. For both tests, a significantly high

Table 7. Kruskal–Wallis H-test between yield distributions grouped by countries and by crop type. \*\*\* =  $p < 0.0001$ .

Evaluation	Country	Crop
p-value	***	***

Table 8. Pairwise post-hoc comparisons of the yield distributions between individual **countries**. Each cell displays the statistical significance level of the difference between two countries based on the Dunn’s test using the Holm–Bonferroni correction. ns = no significance ( $p \geq 0.05$ ), \* =  $p < 0.05$ , \*\* =  $p < 0.01$ , \*\*\* =  $p < 0.001$ , \*\*\*\* =  $p < 0.0001$ .

Comparison	Argentina	Brazil	Germany	Uruguay
Argentina	-	****	****	****
Brazil	****	-	ns	****
Germany	****	ns	-	****
Uruguay	****	****	****	-

Table 9. Pairwise post-hoc comparisons of the yield distributions between individual **crops**. Each cell displays the statistical significance level of the difference between two crops based on the Dunn’s test using the Holm–Bonferroni correction. ns = no significance ( $p \geq 0.05$ ), \* =  $p < 0.05$ , \*\* =  $p < 0.01$ , \*\*\* =  $p < 0.001$ , \*\*\*\* =  $p < 0.0001$ .

Comparison	Corn	Rapeseed	Soybean	Wheat
Corn	-	****	****	****
Rapeseed	****	-	ns	****
Soybean	****	ns	-	****
Wheat	****	****	****	-

Table 10. Kruskal–Wallis H-test between yield distributions for soybean in Argentina, grouped by years and by regions. \*\*\* =  $p < 0.0001$ .

Evaluation	Year	Region
p-value	***	***

Table 11. Pairwise post-hoc comparisons of the yield distributions for soybean in Argentina between individual **years**. Each cell displays the statistical significance level of the difference between two years based on Dunn’s test using the Holm–Bonferroni correction. ns = no significance ( $p \geq 0.05$ ), \* =  $p < 0.05$ , \*\* =  $p < 0.01$ , \*\*\* =  $p < 0.001$ , \*\*\*\* =  $p < 0.0001$ .

Year	2017	2018	2019	2020	2021	2022
2017	-	****	***	****	****	****
2018	****	-	ns	***	****	****
2019	***	ns	-	**	****	****
2020	****	***	**	-	****	****
2021	****	****	****	****	-	****
2022	****	****	****	****	****	-

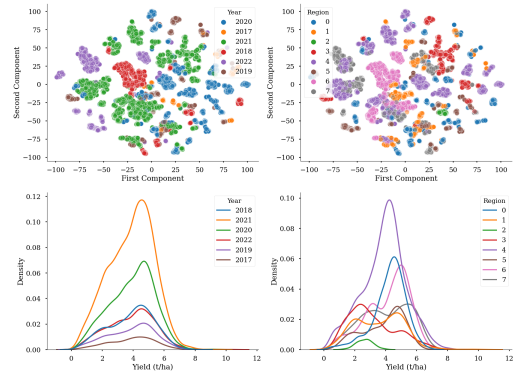


Figure 17. Visualization of the surface reflectance and yield data distribution only for soybean in Argentina, grouped by years and regions. Top left: t-SNE plot of the S2 surface reflectance data grouped by **years**. Top right: t-SNE plot of the S2 surface reflectance data grouped by **region**. Bottom left: Kernel density estimation plot for the target yield data grouped by **years**. Bottom right: Kernel density estimation plot for the target yield data grouped by **regions**.

Table 12. Pairwise post-hoc comparisons of the yield distributions for soybean in Argentina between individual **regions**. Each cell displays the statistical significance level of the difference between two regions based on Dunn’s test using the Holm–Bonferroni correction. ns = no significance ( $p \geq 0.05$ ), \* =  $p < 0.05$ , \*\* =  $p < 0.01$ , \*\*\* =  $p < 0.001$ , \*\*\*\* =  $p < 0.0001$ .

Region	0	1	2	3	4	5	6	7
0	-	****	****	****	****	****	**	****
1	****	-	****	****	****	****	****	****
2	****	****	-	**	****	****	****	****
3	****	****	**	-	****	****	****	****
4	****	****	****	****	-	ns	****	ns
5	****	****	****	****	ns	-	ns	ns
6	**	****	****	****	****	ns	-	***
7	****	****	****	****	ns	ns	***	-

p-value is reported ( $p < 0.0001$ ). This means that at least one year and one region have a significantly different yield distribution compared to the rest of the groups. To compare each pair of yield distributions, a post hoc test is performed using Dunn’s test. The pairwise comparison between the years is depicted in Tab. 11. We highlight that most pairwise comparisons indicate significantly different distributions. Only between 2018 and 2019 is there no significance. The pairwise comparison between the regions is depicted in Tab. 12. As in previous years, the data distribution across regions is mostly significantly different. However, individual regions do not show a significantly different distribution. For instance, region 5 shows no significant difference between regions 4, 6, and 7.

In conclusion, the results undermine the assumption that the distribution across years and regions is mostly significantly different, which increases the difficulty of generalization.

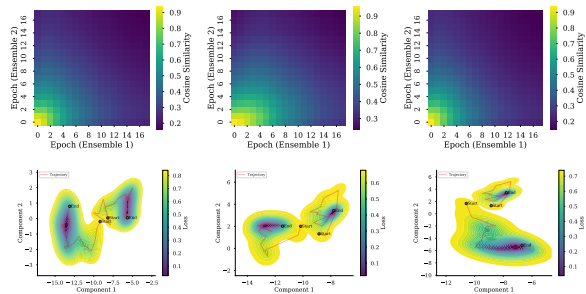


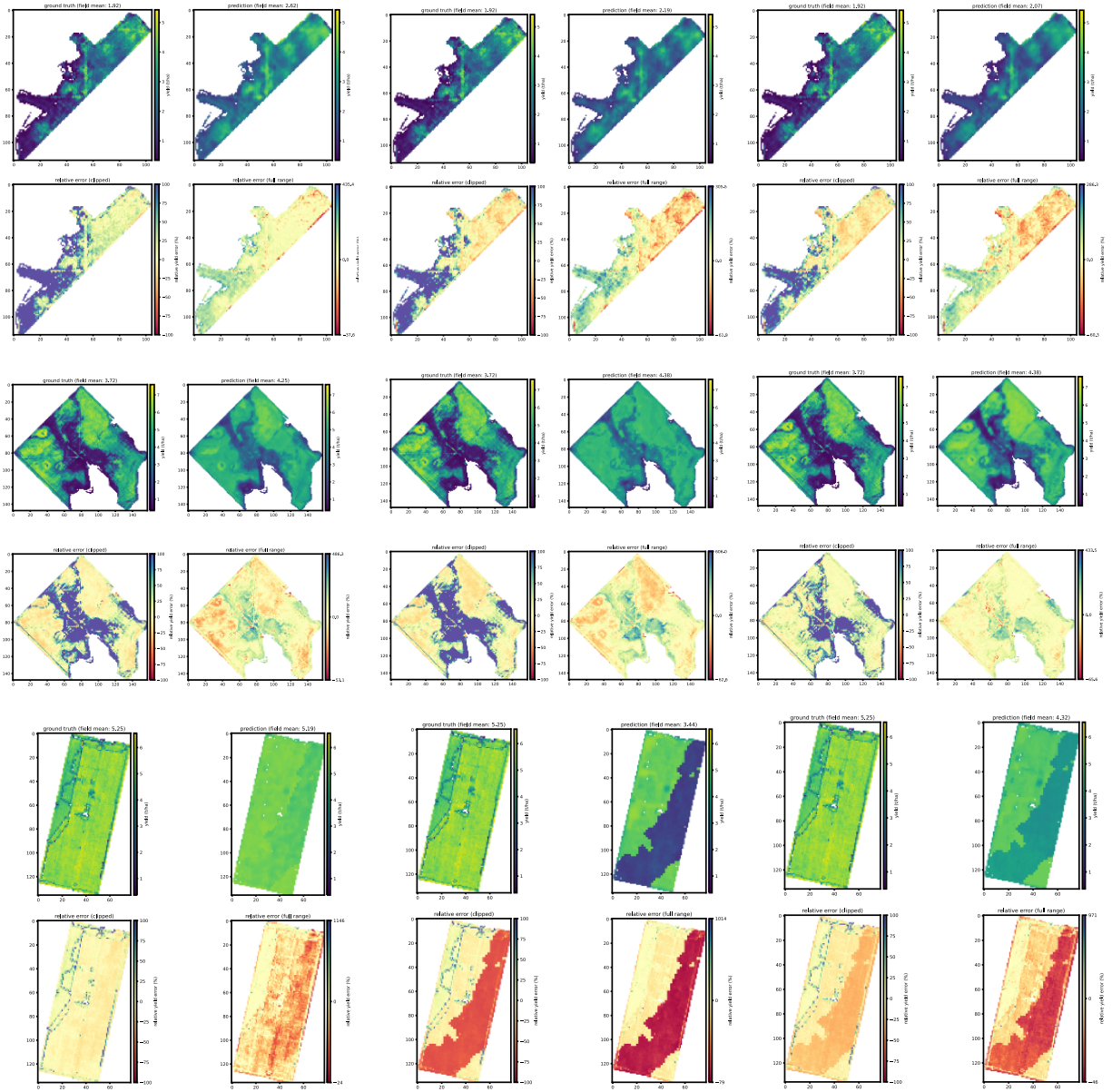
Figure 18. Visualization of the weight space diversity during model training. Top: Cosine similarity between pairs of ensemble members during training. Bottom: PCA plot of the weight space during training, together with the loss. The trajectory in the weight space is highlighted in red from start to end.

### A.3. Further Experiments

Fig. 19 shows additional example predictions for single fields (image), illustrating the ground-truth yield map, the predicted yield map, and the pixel-wise error (clipped and full range). We show the results for the standard 10-fold cv experiment and the temporal (LOYO) and spatial (LORO) experiments. Note that for some fields, the LOYO and LORO experiments exhibit severe performance collapse, as evidenced by high pixel-wise errors and areas of mode collapse (only single-scalar predictions with no spatial variability).

In Tab. 13 we provide a full replication for the 10-fold cv baseline experiment of all benchmark models for each dataset at the field level and at the subfield (pixel)-level in Tab. 13. We further provide a full replication of all benchmark models on the OOD experiments, namely the temporal (LOYO) and spatial (LORO) splitting. For the temporal splitting, the field level results are given in Tab. 15 and for the subfield (pixel)-level in Tab. 16. For the spatial splitting (LORO), the field level results are given in Tab. 17 and for the subfield (pixel)-level in Tab. 18.

Fig. 18 provides more examples of the cosine similarity in weight space over the training epochs, together with the trajectory in weight space.



(a) Standard 10-Fold CV

(b) LOYO CV

(c) LOFO CV

Figure 19. Example ground truth and predicted yield maps together with the pixel-wise error for the standard CV, the LORO, and LOYO CV under distribution shift. The results are generated with the LSTM model [36], trained on S2 data.

Table 13. Results for the RMSE (t/ha) ( $\downarrow$ ) and the  $R^2$ -score ( $\uparrow$ ) for different models and datasets for the standard 10-fold cross-validation at the field level. ARG = Argentina, BRA = Brazil, GER = Germany, URG = Uruguay. C = corn, R = rapeseed, S = soybean, W = wheat.

Evaluation				Field-Level																		
Splitting	Modalities	Fusion Method	Model	ARG-C		ARG-S		ARG-W		BRA-C		BRA-S		GER-R		GER-W		URG-S				
				$R^2$	RMSE	$R^2$	RMSE	$R^2$	RMSE	$R^2$	RMSE	$R^2$	RMSE	$R^2$	RMSE	$R^2$	RMSE	$R^2$	RMSE	$R^2$		
CV10	Sentinel-2	X	3D-ConvLSTM	0.84	1.13	0.79	0.55	0.92	0.62	0.82	0.74	0.76	0.34	0.73	0.42	0.81	0.58	0.65	1.12	0.77	0.51	
			3D-LSTM	0.74	1.44	0.77	0.58	0.89	0.70	0.82	0.74	0.76	0.34	0.63	0.48	0.82	0.57	0.54	1.28	0.73	0.56	
			LSTM	0.79	1.29	0.72	0.64	0.85	0.84	0.75	0.88	0.72	0.68	0.39	0.62	0.49	0.62	0.83	0.56	2.60	0.66	0.62
			Transformer	0.77	1.35	0.73	0.63	0.87	0.78	0.79	0.82	0.75	0.35	0.62	0.49	0.75	0.67	0.67	0.56	1.26	0.72	0.56
			AF	0.84	1.12	0.84	0.84	0.92	0.86	0.82	0.84	0.70	0.81	0.31	0.60	0.71	0.43	0.80	0.60	0.81	0.96	0.81
	Sentinel-2 + ADM	Feature Fusion	MMGF	0.79	1.30	0.82	0.82	0.89	0.69	0.72	0.76	0.86	0.31	0.59	0.51	0.75	0.68	0.77	0.90	0.75	0.53	
			3D-ConvLSTM	0.12	2.63	0.82	0.52	0.89	0.72	0.83	0.83	0.72	0.81	0.34	0.66	0.42	0.78	0.63	0.70	1.03	0.81	0.50
			3D-LSTM	0.40	1.56	0.76	0.59	0.84	0.87	0.77	0.82	0.84	0.71	0.35	0.66	0.72	0.42	0.81	0.67	1.17	0.77	0.52
			LSTM	0.76	1.39	0.72	0.64	0.85	0.83	0.81	0.81	0.78	0.64	0.42	0.60	0.52	0.81	0.58	0.63	1.16	0.72	0.56
			Transformer	0.69	1.55	0.72	0.64	0.83	0.87	0.79	0.81	0.71	0.66	0.37	0.60	0.61	0.49	0.76	0.65	0.61	1.19	0.73

Table 14. Results for the RMSE (t/ha) ( $\downarrow$ ) and the  $R^2$ -score ( $\uparrow$ ) for different models and datasets for the standard 10-fold cross-validation at the subfield (pixel)-level. ARG = Argentina, URG = Uruguay, GER = Germany, BRA = Brazil. S = soybean, R = rapeseed, W = wheat, C = Corn.

Evaluation				Subfield (Pixel)-Level																				
Splitting	Modalities	Fusion Method	Model	ARG-C		ARG-S		ARG-W		BRA-C		BRA-S		BRA-W		GER-R		GER-W		URG-S				
				$R^2$	RMSE	$R^2$	RMSE	$R^2$	RMSE	$R^2$	RMSE	$R^2$	RMSE	$R^2$	RMSE	$R^2$	RMSE	$R^2$	RMSE	$R^2$	RMSE	$R^2$		
CV10	Sentinel-2	X	3D-ConvLSTM	0.65	2.2	0.65	0.9	0.79	1.1	0.45	2.13	0.39	0.93	0.24	1.37	0.48	1.2	0.3	2.46	0.39	1.23	0.51		
			3D-LSTM	0.59	2.37	0.65	0.9	0.79	1.08	0.46	2.13	0.37	0.93	0.24	1.37	0.48	1.2	0.3	2.46	0.39	1.23	0.51		
			LSTM	0.59	2.4	0.65	0.96	0.74	1.21	0.42	2.12	0.34	0.98	0.22	1.39	0.36	1.33	-0.32	3.38	0.37	1.26	0.51		
			Transformer	0.61	2.33	0.62	0.94	0.76	1.17	0.44	2.16	0.38	0.94	0.22	1.39	0.44	1.25	0.32	3.42	0.43	1.24	0.51		
	Sentinel-2 + ADM	Input Fusion	X	3D-ConvLSTM	0.7	2.03	0.73	0.98	0.79	1.04	0.84	0.96	0.67	0.46	2.12	0.44	0.9	0.24	1.37	0.48	1.20	0.51	0.44	
				3D-LSTM	0.65	2.2	0.73	0.98	0.84	1.07	0.42	2.19	0.42	0.92	0.2	1.4	0.44	1.26	0.51	2.2	0.36	1.19	0.44	
				LSTM	0.08	3.58	0.68	0.87	0.78	1.11	0.46	2.12	0.38	0.94	0.23	1.38	0.42	1.27	0.51	0.39	2.3	0.41	0.61	0.44
				Transformer	0.36	2.97	0.64	0.91	0.78	1.13	0.43	2.12	0.41	0.93	0.24	1.37	0.49	1.2	0.37	3.34	0.4	1.22	0.51	
	Sentinel-2 + ADM	Input Fusion	X	3D-ConvLSTM	0.58	2.41	0.59	0.98	0.75	1.2	0.43	2.18	0.37	0.99	0.21	1.4	0.47	1.21	0.51	0.35	2.38	0.39	1.23	0.51
				3D-LSTM	0.52	2.57	0.58	0.98	0.7	1.3	0.44	2.17	0.37	0.96	0.22	1.39	0.45	1.24	0.38	2.32	0.39	1.23	0.51	
				LSTM	0.52	2.57	0.58	0.98	0.7	1.3	0.44	2.17	0.37	0.96	0.22	1.39	0.45	1.24	0.38	2.32	0.39	1.23	0.51	
				Transformer	0.52	2.57	0.58	0.98	0.7	1.3	0.44	2.17	0.37	0.96	0.22	1.39	0.45	1.24	0.38	2.32	0.39	1.23	0.51	

Table 15. Results for the RMSE (t/ha) ( $\downarrow$ ) and the  $R^2$ -score ( $\uparrow$ ) for different models and datasets for the Leave-One-Year-Out cross-validation at the field level. ARG = Argentina, BRA = Brazil, GER = Germany, URG = Uruguay. C = corn, R = rapeseed, S = soybean, W = wheat.

Evaluation				Field-Level																				
Splitting	Modalities	Fusion Method	Model	ARG-C		ARG-S		ARG-W		BRA-C		BRA-S		BRA-W		GER-R		GER-W		URG-S				
				$R^2$	RMSE	$R^2$	RMSE	$R^2$	RMSE	$R^2$	RMSE	$R^2$	RMSE	$R^2$	RMSE	$R^2$	RMSE	$R^2$	RMSE	$R^2$	RMSE	$R^2$		
LOYO	Sentinel-2	X	3D-ConvLSTM	0.61	1.75	0.65	0.71	0.78	1.00	0.20	1.58	0.57	0.55	-0.06	0.82	0.43	1.02	0.23	1.60	0.63	0.55	0.67		
			3D-LSTM	0.53	1.93	0.66	0.70	0.77	1.02	0.36	1.42	0.12	0.65	-0.33	0.92	0.41	0.53	0.00	1.97	0.62	0.66	0.67		
			LSTM	0.46	2.06	0.50	0.82	0.77	1.10	0.29	1.59	0.06	0.61	0.67	-0.05	1.00	0.82	-0.07	1.39	0.66	0.61	0.67		
			Transformer	0.57	1.83	0.52	0.84	0.78	1.04	0.39	1.19	0.59	0.28	0.59	-0.09	0.83	0.33	1.10	0.56	1.76	0.60	0.68	0.66	
	Sentinel-2 + ADM	Input Fusion	X	3D-ConvLSTM	0.33	2.90	0.77	0.59	0.81	0.93	0.42	1.57	0.38	0.54	0.14	0.74	0.44	0.90	0.19	1.71	0.56	0.53	0.73	
				3D-LSTM	0.31	2.33	0.66	0.70	0.36	1.71	0.39	1.38	0.38	0.55	0.08	0.83	0.15	1.24	-0.03	1.92	0.48	0.77	0.67	
				LSTM	-0.01	2.82	0.67	0.69	0.40	1.03	0.45	1.32	0.34	0.56	-0.15	0.86	0.49	0.56	0.96	0.37	1.39	0.56	0.71	0.67
				Transformer	0.18	2.54	0.62	0.74	0.78	0.99	0.28	1.50	0.45	0.52	-0.07	0.82	0.26	0.88	1.15	0.59	1.75	0.50	0.74	0.67
	Sentinel-2 + ADM	Input Fusion	X	3D-ConvLSTM	0.42	2.13	0.62	0.74	0.74	1.08	0.29	1.49	0.26	0.59	0.12	0.75	0.58	0.87	0.33	1.56	0.54	0.50	0.73	
				3D-LSTM	0.33	2.30	0.65	0.72	0.75	1.08	0.25	1.54	0.20	0.59	0.11	0.75	0.58	0.87	0.33	1.56	0.54	0.50	0.73	
				LSTM	0.33	2.30	0.65	0.72	0.75	1.08	0.25	1.54	0.20	0.59	0.11	0.75	0.58	0.87	0.33	1.56	0.54	0.50	0.73	
				Transformer	0.33	2.30	0.65	0.72	0.75	1.08	0.25	1.54	0.20	0.59	0.11	0.75	0.58	0.87	0.33	1.56	0.54	0.50	0.73	

Table 16. Results for the RMSE (t/ha) ( $\downarrow$ ) and the  $R^2$ -score ( $\uparrow$ ) for different models and datasets for the Leave-One-Region-Out (LORO) cross-validation at the field level. ARG = Argentina, BRA = Brazil, GER = Germany, URG = Uruguay. C = corn, R = rapeseed, S = soybean, W = wheat.

Evaluation				Subfield (Pixel)-Level																			
Splitting	Modalities	Fusion Method	Model	ARG-C		ARG-S		ARG-W		BRA-C		BRA-S		BRA-W		GER-R		GER-W		URG-S			
				$R^2$	RMSE	$R^2$	RMSE	$R^2$	RMSE	$R^2$	RMSE	$R^2$	RMSE	$R^2$	RMSE	$R^2$	RMSE	$R^2$	RMSE	$R^2$	RMSE	$R^2$	
LOYO	Sentinel-2	X	3D-ConvLSTM	0.48	2.68	0.54	1.04	0.67	1.38	0.15	2.67	0.19	1.08	0.06	1.52	0.21	1.49	0.53	2.71	0.39	0.34	1.28	
			3D-LSTM	0.47	2.70	0.55	1.03	0.68	1.35	0.21	2.56	0.13	1.12	0.05	1.54	0.18	1.46	0.51	2.93	0.39	0.35	1.27	
			LSTM	0.39	2.90	0.46	1.12	0.61	1.49	0.15	2.66	0.04	1.18	0.08	1.57	0.15	1.71	-0.08	3.82	0.33	0.33	1.29	
			Transformer	0.42	2.83	0.48	1.10	0.68	1.38	0.20	2.60	0.16	1.10	0.05	1.53	0.14	1.49	0.55	3.13	0.33	0.33	1.29	
	Sentinel-2 + ADM	Input Fusion	X	3D-ConvLSTM	0.35	2.99	0.66	0.89	0.72	1.27	0.23	2.43	0.29	1.02	0.11	1.48	0.25	1.40	0.62	2.70	0.32	0.32	1.30
				3D-LSTM	0.30	3.10	0.57	1.00	0.68	2.03	0.26	2.53	0.24	1.05	0.06	1.52	0.28	1.63	0.60	2.86	0.37	0.28	1.31
				LSTM	-0.02	3.75	0.56	1.01	0.64	1.43	0.26	2.47	0.19	1.08	0.04	1.54	0.22	1.47	0.52	2.60	0.31	0.31	1.31
				Transformer	0.25	3.23	0.53	1.05	0.68	1.34	0.22	2.55	0.22	1.06	0.06	1.52	0.13	1.56	0.53	2.72	0.30	0.30	1.32
	Sentinel-2 + ADM	Input Fusion	X	3D-ConvLSTM	0.44	2.78	0.49	1.09	0.61	1.49	0.07	0.19	0.60	0.60	1.17	0.09	0.99	0.31	1.00	0.88	2.67	0.30	0.30
				3D-LSTM	0.44	2.78	0.49	1.09	0.61	1.49	0.07	0.19	0.60	0.60	1.17	0.09	0.99	0.31	1.00	0.88	2.67	0.30	0.30
				LSTM	0.24	3.24	0.49	1.09	0.61	1.49	0.07	0.19	0.60	0.60	1.17	0.09	0.99	0.31	1.00	0.88	2.67	0.30	0.30
				Transformer	0.24	3.24	0.49	1.09	0.61	1.49	0.07	0.19	0.60	0.60	1.17	0.09	0.99	0.31	1.00	0.88	2.67	0.30	0.30

Table 17. Results for the RMSE (t/ha) ( $\downarrow$ ) and the  $R^2$ -score ( $\uparrow$ ) for different models and datasets for the Leave-One-Region-Out (LORO) cross-validation at the field level. ARG = Argentina, BRA = Brazil, GER = Germany, URG = Uruguay. C = corn, R = rapeseed, S = soybean, W = wheat.

Evaluation				Field-Level																				
Splitting	Modalities	Fusion Method	Model	ARG-C		ARG-S		ARG-W		BRA-C		BRA-S		BRA-W		GER-R		GER-W		URG-S				
				$R^2$	RMSE	$R^2$	RMSE	$R^2$	RMSE	$R^2$	RMSE	$R^2$	RMSE	$R^2$	RMSE	$R^2$	RMSE	$R^2$	RMSE	$R^2$	RMSE	$R^2$		
LOYO	Sentinel-2	X	3D-ConvLSTM	0.68	1.59	0.69	0.67	0.80	0.95	0.55	1.19	0.37	0.55	0.47	0.58	0.25	1.16	0.14	1.76	0.68	0.50	0.61		
			3D-LSTM	0.49	2.00	0.70	0.67	0.81	0.92	0.59	1.13	0.41	0.53	0.48	0.57	0.26	1.16	-0.08	1.98	0.65	0.60	0.63		
			LSTM	0.47	2.05	0.64	0.72	0.73	1.18	0.39	1.39	0.27	0.59	0.27	0.69	0.68	-0.16	1.45	0.49	2.71	0.36	0.63	0.65	
			Transformer	0.55	1.88	0.59	0.77	0.71	1.16	0.35	1.43	0.30	0.57	0.37	0.63	0.63	0.17	1.22	0.50	1.89	0.63	0.61	0.65	
	Sentinel-2 + ADM	Input Fusion	X	3D-ConvLSTM	0.53	1.92	0.78	0.57	0.87	0.78	0.55	1.19	0.33	0.42	0.52	0.55	0.20	1.29	-0.85	3.21	0.60	0.60	0.68	
				3D-LSTM	0.59	1.80	0.65	0.77	0.78	1.00	0.10	1.68	0.33	0.58	0.23	0.70	-0.56	1.68	0.60	-0.95	2.65	0.54	0.73	0.68
				LSTM	0.68	1.59	0.69	0.67	0.80	0.95	0.55	1.19	0.37	0.55	0.47	0.58	0.25	1.16	0.14	1.76	0.68	0.60	0.61	
				Transformer	0.23	2.46	0.66	0.70	0.74	1.09	0.65	1.39	0.29	0.47	0.51	0.48	-0.16	0.58	0.30	1.32	0.50	0.62	0.66	
	Sentinel-2 + ADM	Input Fusion	X	3D-ConvLSTM	0.43	2.12	0.66	0.70	0.61	1.34	0.38	1.39	0.29	0.58	0.26	0.69	0.20	1.2	-1.39	2.94	0.56	0.61	0.68	
				3D-LSTM	0.44	2.10	0.67	0.70	0.63	1.30	0.52	1.22	0.20	0.62	0.36	0.64	0.20	1.18	1.22	2.20	0.59	0.70	0.51	
				LSTM	0.44	2.10	0.67	0.70	0.63	1.30	0.52	1.22	0.20	0.62	0.36	0.64	0.20	1.18						

## B. Dataset Datasheet

In the following, we will refer to the authors and their institutions in the paper as the “publishing institution.”

### B.1. Motivation

The motivation to create the dataset was to advance large-scale crop yield prediction. Data creation was funded by industry and public stakeholders.

### B.2. Dataset Creation

The dataset can be used for crop yield prediction as described in the main paper. The dataset can be used for further research or other directions, such as crop classification, time series analysis, or land cover classification, provided that data privacy and the data license are complied with.

### B.3. Data Collection Process

All the yield data was collected in collaboration with local farmers, data providers, and industry companies. Each data provider had to sign a data-sharing agreement with the publishing institution. Data providers were compensated for data collection and data sharing on an area basis.

The EO data was acquired by the publishing institution from only publicly available data providers under the licenses of those providers.

The dataset does not contain any other data from other publicly available data sources or libraries and is entirely self-contained. The data contains sensitive information about the geographic locations of the local data providers. Each data point is georeferenced. We provide georeferences to facilitate downstream research, including additional data acquisition activities beyond the current dataset.

The dataset may have errors and inconsistencies arising during the data collection process. This may include errors in the yield data collection and preprocessing processes. The dataset may include local biases to individual regions, crop types, and years. Moreover, the data may include data imbalances that could bias downstream models. We provide detailed metadata to help researchers account for these limitations.

### B.4. Data Preprocessing

The data preprocessing was described in detail in Sec. 4.2. Statistics to calculate the target scaled yield (dry yield) are given in Tab. 6.

### B.5. Data Distribution

The publishing institution owns the data in its entirety and will grant access under a custom license.

### B.6. Data Maintenance

The dataset is entirely self-contained. The publishing institution will maintain the data entirely on its internal servers.

The publishing institution will maintain the data access portal and provide sufficient information to contact it. The data will be continuously maintained and updated. Updates will be documented on the publishing institution’s webpage. We welcome community contributions that improve the data quality.

The dataset is available at <https://yieldsat.github.io/>

### B.7. Legal and Ethical Considerations

The data is owned entirely by the publishing institution. Access will be granted under the described license and data privacy regulation. The data does not contain any information about human subjects or any other confidential information. However, sensitive information is provided in the georeferences of the data.

This dataset is intended for academic research in a non-commercial setting. There are no expectations for this dataset to produce models that can generalize to real-world data, and the data usage is aimed for academic research, for example: to build new yield estimation models, evaluate yield estimation models, and study the yield estimation process itself.