| cases | doc_1 | | doc_2 | | decision | id |
|---|---|---|---|---|---|---|
| | authors | <ul><li>Yue He</li><li>Chen Chen</li><li>Jing Zhang</li><li>Juhua Liu</li><li>Fengxiang He</li><li>Chaoyue Wang</li><li>Bo Du</li></ul> | authors | <ul><li>Yue He</li><li>Chen Chen</li><li>Jing Zhang</li><li>Juhua Liu</li><li>Fengxiang He</li><li>Chaoyue Wang</li><li>Bo Du</li></ul> | DUPLICATES | 29 |
| | title | Visual Semantics Allow for Textual Reasoning Better in Scene Text Recognition | title | Visual Semantics Allow for Textual Reasoning Better in Scene Text Recognition | | |
| | publication_date | 2021-12-24 02:43:42+00:00 | publication_date | 2021-12-24 00:00:00 | | |
| | source | SupportedSources.ARXIV | source | SupportedSources.INTERNET_ARCHIVE | | |
| | journal | None | journal | | | |
| | volume | | volume | | | |
| | doi | | doi | | | |
| | urls | <ul><li>http://arxiv.org/pdf/2112.12916v1</li><li>http://arxiv.org/abs/2112.12916v1</li><li>http://arxiv.org/pdf/2112.12916v1</li></ul> | urls | <ul><li>https://web.archive.org/web/20220105103329/https://arxiv.org/pdf/2112.12916v1.pdf</li></ul> | | |
| | id | id8715457885182187645 | id | id433245745519123933 | | |
| | abstract | Existing Scene Text Recognition (STR) methods typically use a language model to optimize the joint probability of the 1D character sequence predicted by a visual recognition (VR) model, which ignore the 2D spatial context of visual semantics within and between character instances, making them not generalize well to arbitrary shape scene text. To address this issue, we make the first attempt to perform textual reasoning based on visual semantics in this paper. Technically, given the character segmentation maps predicted by a VR model, we construct a subgraph for each instance, where nodes represent the pixels in it and edges are added between nodes based on their spatial similarity. Then, these subgraphs are sequentially connected by their root nodes and merged into a complete graph. Based on this graph, we devise a graph convolutional network for textual reasoning (GTR) by supervising it with a cross-entropy loss. GTR can be easily plugged in representative STR models to improve their performance owing to better textual reasoning. Specifically, we construct our model, namely S-GTR, by paralleling GTR to the language model in a segmentation-based STR baseline, which can effectively exploit the visual-linguistic complementarity via mutual learning. S-GTR sets new state-of-the-art on six challenging STR benchmarks and generalizes well to multi-linguistic datasets. Code is available at https://github.com/adeline-cs/GTR. | abstract | Existing Scene Text Recognition (STR) methods typically use a language model to optimize the joint probability of the 1D character sequence predicted by a visual recognition (VR) model, which ignore the 2D spatial context of visual semantics within and between character instances, making them not generalize well to arbitrary shape scene text. To address this issue, we make the first attempt to perform textual reasoning based on visual semantics in this paper. Technically, given the character segmentation maps predicted by a VR model, we construct a subgraph for each instance, where nodes represent the pixels in it and edges are added between nodes based on their spatial similarity. Then, these subgraphs are sequentially connected by their root nodes and merged into a complete graph. Based on this graph, we devise a graph convolutional network for textual reasoning (GTR) by supervising it with a cross-entropy loss. GTR can be easily plugged in representative STR models to improve their performance owing to better textual reasoning. Specifically, we construct our model, namely S-GTR, by paralleling GTR to the language model in a segmentation-based STR baseline, which can effectively exploit the visual-linguistic complementarity via mutual learning. S-GTR sets new state-of-the-art on six challenging STR benchmarks and generalizes well to multi-linguistic datasets. Code is available at https://github.com/adeline-cs/GTR. | | |
| | versions | | versions | | | |