

: ---.

इंटरनेशनल साइंटिफिक जर्नल ऑफ इंजीनियरिंग एंड मैनेजमेंट () : 2583-

वॉल्यूम: 04 अंक: 05 | मई -

: 10.55041/03470

एक अंतरराष्ट्रीय वदिवतापूर्ण | | बहु-वषियक | | ओपन एक्सेस | | सभी प्रमुख डेटाबेस और मेटाडेटा में अनुक्रमण

2025, (सभी अधिकार सुरक्षित) |.. | पेज

एआई-आधारित पीडीएफ अनुवादक

आर। राजेश*1, पी। अशोक*2, एम। साई कृष्णा*

* इंजीनियरिंग कॉलेज, भारत के () विभाग के 1 सहायक प्रोफेसर।

* इंजीनियरिंग कॉलेज, भारत के () विभाग के 2,3 छात्र।

अमूर्त

स्वचालित दस्तावेज़ अनुवाद भाषा की बाधाओं पर काबू पाने और सीमलेस को सुविधाजनक बनाने में महत्वपूर्ण भूमिका निभाता है

वैश्विक उद्योगों में संचार। यह परियोजना प्राकृतिक भाषा प्रसंस्करण (एनएलपी) की शक्ति का उपयोग करती है और ऑप्टिकल चरित्र मान्यता () पीडीएफ दस्तावेजों से पाठ को निकालने, अनुवाद करने और पुनर्निर्माण करने के लिए

उनके मूल लेआउट और स्वरूपण को संरक्षित करना। ट्रांसफार्मर-आधारित मॉडल जैसे जीपीटी और का उपयोग करके

एपीआई का अनुवाद करें, मजबूत पाठ निष्कर्षण उपकरणों के साथ, सिस्टम सटीक और कुशल बहुभाषी प्रदान करता है

अनुवाद। कार्यप्रणाली में , , -, और सहित पायथन पुस्तकालयों को शामिल किया गया है

पाठ मान्यता, अनुवाद और सुधार प्रक्रियाओं का प्रबंधन करने के लिए। यह एआई-संचालित समाधान उद्देश्य है

पहुंच बढ़ाने के लिए, वैश्विक सहयोग को बढ़ावा दें, और बहुभाषी दस्तावेज़ वर्कफ़्लो को सुव्यवस्थित करें विविध क्षेत्र।

कीवर्ड: पीडीएफ अनुवाद, प्राकृतिक भाषा प्रसंस्करण (एनएलपी), ऑप्टिकल चरित्र मान्यता (ओसीआर), ट्रांसफार्मर मॉडल, बहुभाषी दस्तावेज़ प्रसंस्करण।

मैं।

परिचय

1.1. और प्रेरणा:

एक तेजी से वैश्विक और परस्पर जुड़ी दुनिया में, संगठनों को संवाद और वनिमिय करने की आवश्यकता होती है भाषाई सीमाओं के पार जानकारी। पोर्टेबल डॉक्यूमेंट फॉर्मेट (पीडीएफ) में दस्तावेज़ व्यापक रूप से उपयोग किए जाते हैं

स्वास्थ्य सेवा, वित्त, कानून, शिक्षाविद और सरकार जैसे उद्योग उनकी पोर्टेबिलिटी, सुसंगत स्वरूपण, के कारण,

और बहु-प्लेटफॉर्म संगतता। हालांकि, विभिन्न भाषाओं में पीडीएफ दस्तावेजों का अनुवाद करना कई प्रस्तुत करता है

चुनौतियां - विशेष रूप से जब इन दस्तावेजों में जटिल लेआउट, स्कैन की गई छवियां और बहुभाषी सामग्री होती है।

पारंपरिक अनुवाद वधियाँ अक्सर मैनुअल, समय लेने वाली, महंगी और त्रुटियों के लिए प्रवण होती हैं, जिससे वे बनते हैं बड़े पैमाने पर दस्तावेज़ प्रसंस्करण के लिए अक्षम।

1.2 परिचय

आज की तेजी से वैश्वीकृत दुनिया में, वभिन्न में जानकारी संवाद करने और साझा करने की क्षमता भाषाएँ महत्वपूर्ण हैं। संगठन अक्सर विविध हतिधारकों, भागीदारों और ग्राहकों के साथ काम करते हैं जो वभिन्न बोलते हैं

भाषाएं, कुशल, सटीक और स्केलेबल अनुवाद समाधानों के लिए बढ़ती मांग पैदा करती हैं। परंपरागत दस्तावेज़ अनुवाद के तरीके समय लेने वाले, महंगे और मानवीय त्रुटि के लिए प्रवण हैं-वर्षों से व्यवहार करते समय

पीडीएफ जैसे जटिल प्रारूपों में सामग्री के बड़े संस्करणों के साथ।

पीडीएफ दस्तावेजों का उपयोग उनके निश्चित स्वरूपण और क्रॉस-प्लेटफॉर्म संगतता के कारण उद्योगों में व्यापक रूप से किया जाता है।

हालांकि, उनकी कठोर संरचना स्वचालित अनुवाद के लिए चुनौतियों का सामना करती है, खासकर जब वे स्कैन किए गए होते हैं

छवियाँ, टेबल, या बहुभाषी सामग्री। यह बुद्धिमान प्रौद्योगिकियों को अपनाने के लिए आवश्यक बनाता है जो सटीक रूप से कर सकते हैं

ऐसे प्रारूपों से पाठ्य डेटा निकालें, व्याख्या और पुनः इकट्ठा करें।

यह परियोजना स्वचालित पीडीएफ अनुवाद के लिए एआई-संचालित प्रणाली विकसित करके इन चुनौतियों को संबोधित करती है। द्वारा

लीवरेजिंग नेचुरल लैंग्वेज प्रोसेसिंग (एनएलपी), ऑप्टिकल कैरेक्टर रिकग्निशन (ओसीआर), और एडवांस्ड ट्रांसफार्मर-

और अनुवाद जैसे आधारित मॉडल, सिस्टम देशी और स्कैन किए गए पाठ दोनों का अनुवाद कर सकता है जबकि

इंटरनेशनल साइटफिकि जर्नल ऑफ़ इंजीनियरिंग एंड मैनेजमेंट () : 2583-

वॉल्यूम: 04 अंक: 05 | मई -

: 10.55041/03470

एक अंतरराष्ट्रीय वद्वितापूर्ण || बहु-वर्षिक || ओपन एक्सेस || सभी प्रमुख डेटाबेस और मेटाडेटा में अनुक्रमण

2025, (सभी अधिकार सुरक्षित) |.. | पेज

मूल लेआउट को बनाए रखना।, , -, और जैसे पायथन

लाइब्रेरी

एपीआई इस समाधान की बैकबोन बनाता है, जो निष्कर्षण, अनुवाद और पुनर्निर्माण के सहज एकीकरण को सक्षम करता है

वर्कफ्लो।

इस परियोजना का अंतिम लक्ष्य भाषा की बाधाओं को तोड़ना, पहुंच में सुधार करना और वैश्विक बढ़ाना है बुद्धिमान बहुभाषी दस्तावेज़ प्रसंस्करण के माध्यम से सहयोग।

। साहित्यिक निरीक्षण

स्वचालित पीडीएफ अनुवाद प्रणालियों का विकास ऑप्टिकल चरित्र मान्यता में प्रगतिको एकीकृत करता है (), प्राकृतिक भाषा प्रसंस्करण (), और ट्रांसफार्मर-आधारित मॉडल। यह सर्वेक्षण प्रमुख

अनुसंधान पर प्रकाश डालता है

योगदान जो ऐसी प्रणालियों के डिजाइन और कार्यान्वयन को सूचित करते हैं।

2.1। पाठ निष्कर्षण के लिए इंजन

स्मृति ने टेसरेक्ट को पेश किया, जो एक ओपन-सोर्स ओसीआर इंजन है जो वभिन्न भाषाओं में पाठ को पहचानने

में सकृषम है और

सक्रुपिट। की अनुकूलनशीलता और सटीकता इसे स्कैन कए गए दस्तावेजों से पाठ निकालने के लिए एक मूल्यवान उपकरण बनाती है, ए
स्वचालति पीडीएफ अनुवाद वर्कफ्लोज में महत्वपूर्ण कदम।

2.2। संयुक्त रूप से संरेखति करने और अनुवाद करने के लिए संयुक्त रूप से सीखने से तंत्रिका मशीन अनुवाद बहडानाऊ एट अल। एक प्रारंभिक मॉडल प्रस्तुत कया जो संयुक्त रूप से संरेखति करना और अनुवाद करना सीखता है, सीमाओं को संबोधति करता है
पारंपरिक अनुक्रम-से-अनुक्रम मॉडल। उनका दृष्टिकोण एक ध्यान तंत्र को एकीकृत करता है जो मॉडल को सकृषम करता है
अनुवाद के दौरान स्रोत वाक्य के प्रासंगिक भागों पर ध्यान केंद्रति करने के लिए, लंबे अनुक्रमों पर प्रदर्शन में सुधार।

2.3। ध्यान-आधारति तंत्रिका मशीन अनुवाद के लिए प्रभावी दृष्टिकोण
लुओग एट अल। वैश्विक और स्थानीय सहति तंत्रिका मशीन अनुवाद में वभिन्नि ध्यान तंत्रों का पता लगाया ध्यान दे मॉडल। उनके नष्टिकरणों से संकेत मलिता है कधियान-आधारति मॉडल पारंपरिक रूप से बेहतर प्रदर्शन करते हैं
अनुवाद प्रणाली, वशिष रूप से लंबे वाक्यों और जटलि संरचनाओं को संभालने में।

2.4। ध्यान आप सभी की जरूरत है (ट्रांसफार्मर वास्तुकला)
वासवानी एट अल। द्वारा फाउंडेशनल वर्क, शीर्षक "अटेंशन इज़ ऑल यू नीड," ने ट्रांसफार्मर का परिचय दिया आर्कटिकचर, क्रांति मशीन अनुवाद और एनएलपी कार्य। स्व-कृषण तंत्र पर मॉडल की निर्भरता समानांतर प्रसंस्करण के लिए अनुमति देता है और लंबी दूरी की निर्भरता के हैंडलिंग में सुधार करता है, नए बेचमार्क सेट करता है
अनुवाद की गुणवत्ता।

2.5। की तंत्रिका मशीन अनुवाद प्रणाली
वू एट अल। एक तंत्रिका मशीन अनुवाद प्रणाली के लिए के संक्रमण को वसितृत करे जो गहरे नेटवर्क का उपयोग करता है
ध्यान तंत्र। इस प्रणाली ने अनुवाद सटीकता और प्रवाह में पर्याप्त सुधार प्राप्त कया, बड़े पैमाने पर अनुप्रयोगों में तंत्रिका दृष्टिकोण की व्यावहारिक व्यवहार्यता का प्रदर्शन।

2.6। बर्ट: भाषा समझ के लिए गहरी द्वदिशि ट्रांसफार्मर का पूर्व-प्रशिक्षण
डेवलनि एट अल। बर्ट का परिचय, एक पूर्व-प्रशिक्षति गहरी द्वदिशि ट्रांसफार्मर मॉडल जसिने अतृप्त है-वभिन्नि एनएलपी कार्यों में कला परिणाम। बर्ट की वास्तुकला दोनों दिशाओं में संदर्भ की गहरी समझ के लिए अनुमति देती है,
प्रश्न उत्तर और भाषा के अनुमान जैसे कार्यों के लिए इसे अत्यधिक प्रभावी बनाना।

2.7। तंत्रिका अनुवाद मॉडल में व्याख्यात्मक ध्यान तंत्र
जेनकेल एट अल। ट्रांसफार्मर आर्कटिकचर के लिए एक एक्सटेंशन पेश कया जो शामिल करके शब्द संरेखण को बढ़ाता है
व्याख्यात्मक ध्यान तंत्र। उनका मॉडल एनकोडर जानकारी पर ध्यान केंद्रति करता है, अधिक सुविधा प्रदान करता है
प्रशिक्षण के दौरान शब्द-संरेखण डेटा की आवश्यकता के बिना सटीक संरेखण। यह दृष्टिकोण काफी है ट्रांसफार्मर ध्यान सकृषिणों की भोली व्याख्याओं को कम करता है और संरेखण गुणवत्ता को प्राप्त करता है ++ जैसे पारंपरिक उपकरणों के लिए।

2.8।ट्रांसफार्मर के लिए व्याख्या योग्य शब्द संरेखण
डौ और न्युबगि ने व्याख्याता को बढ़ाने के लिए हेड प्रूनगि और पर्यवेक्षण संरेखण तकनीकों की शुरुआत की

इंटरनेशनल साइंटिफिक जर्नल ऑफ इंजीनियरिंग एंड मैनेजमेंट () : 2583-

वॉल्यूम: 04 अंक: 05 | मई -

: 10.55041/03470

एक अंतरराष्ट्रीय वद्वितापूर्ण || बहु -वर्षिक || ओपन एक्सेस || सभी प्रमुख डेटाबेस और मेटाडेटा में
अनुक्रमण

2025, (सभी अधिकार सुरक्षित) |.. | पेज

ट्रांसफार्मर में ध्यान तंत्र।उनका दृष्टिकोण कुंजी का चयन करके अधिक सटीक शब्द संरेखण देता है
ध्यान दे, हालांकि इसके लिए अतिरिक्त पर्यवेक्षण और संरेखण डेटा की आवश्यकता होती है।

2.9।ढाल-आधारित स्व-संयोग मानचित्र विश्लेषण

बार्कन एट अल।प्रस्तावित ग्रेड-एसएम, ट्रांसफार्मर भविष्यवाणियों को स्वयं के माध्यम से समझने के लिए
एक ढाल-आधारित विधि

ध्यान दें।ग्रेड-एसएम इनपुट तत्वों की पहचान करता है जो मॉडल निर्णयों को सबसे महत्वपूर्ण रूप से प्रभावित
करते हैं, प्रदान करते हैं

ट्रांसफार्मर-आधारित भाषा मॉडल के आंतरिक कामकाज में अंतर्दृष्टि।भूलांकन से पता चलता है कि ग्रेड-सैम
मौजूदा व्याख्या तकनीकों पर पर्याप्त सुधार प्रदान करता है।

2.10।बड़े भाषा मॉडल में अनुवाद तंत्र की खोज

झांग एट अल।बड़े भाषा मॉडल (एलएलएम) के भीतर अनुवाद तंत्र को समझने के लिए एक अध्ययन किया।

पथ पैचिंग तकनीकों का उपयोग करते हुए, उन्होंने पाया कि ध्यान देने वाले सरि का एक वरिल सबसेट (5%से कम)

मुख्य रूप से अनुवाद कार्यों की सुविधा प्रदान करता है।केवल इन विशेष सरि को ठीक करने के द्वारा, उन्होंने

अनुवाद प्राप्त किया

पूर्ण-पैरामीटर ट्यूनिंग की तुलना में सुधार, लक्षित मॉडल समायोजन की दक्षता को उजागर करते हुए।

2.11।तुलना तालिका: ध्यान-आधारित मॉडल पर साहित्य समीक्षा

नहीं।

कागज का शीर्षक /

केंद्र

लेखक

वर्ष

कार्यप्रणाली प्रमुख निष्कर्ष

सीमाएँ

/

टिप्पणी

टेसरेक्ट

ओसीआर इंजन

लोहार

एलएसटीएम आधारित

ओसीआर

साथ

भाषा

सहायता
शुद्ध
बहुभाषी
मूलपाठ
मान्यता
छवियों से
संघर्ष
साथ
अत्यधिक शोर या
वर्षा
दस्तावेज़

तंत्रिका
मशीन
द्वारा अनुवाद करना
संयुक्त रूप से
सीखना
को
संरक्षित
और
अनुवाद
बहडानौ
एट अल।

एनकोडर-
के साथ डिकोडर
ध्यान
तंत्र
उन्नत
का अनुवाद
लंबा
दृश्यों
कम्प्यूटेशनल के रूप में
गैर से धीमा
ध्यान
2 मॉडल

असरदार
के लिए दृष्टिकोण
ध्यान-
आधारित तंत्रिका
मशीन
अनुवाद

एट
अल।

वैश्वकि
बनाम
स्थानीय
ध्यान
में
2
मॉडल
स्थानीय
ध्यान
तेजी से मॉडल
और
अधिक
शुद्ध
में
नश्चिति
संदर्भों
आरएनएन तक सीमति
चौखटे

:

तंत्रिका
मशीन
अनुवाद
प्रणाली
वू एट अल।

गहरी
ध्यान के साथ
के लिए
अंत है-
अंत एमटी
प्रमुख
सुधार
प्रवाह में और
अनुवाद
पैमाने पर गुणवत्ता
उच्च
वलिंब
बनिा
अनुकूलन

ध्यान
है
आप सभी की जरूरत
वासवानी एट
अल।

ट्रांसफार्मर
वास्तुकला
का उपयोग करते हुए
खुद-
ध्यान
पेश किया
ट्रांसफार्मर;
बेहतर प्रदर्शन किया

मे
अनुवाद
कार्य
आवश्यक है
बड़ा
डेटासेट
और
गणना
संसाधन

बर्टः
द्विदिशि
ट्रान्सफॉर्मर
डेवलपि
एट
अल।

पूर्व प्रशिक्षण
नकाबपोश के साथ
भाषा
मॉडलिंग
मज़बूत
प्रासंगिक
समझ;
नहीं
अनुकूलति
के लिए
अनुवाद
कार्य
वर्षीय रूप से

इंटरनेशनल साइंटिफिक जर्नल ऑफ़ इंजीनियरिंग एंड मैनेजमेंट () : 2583-
वॉल्यूम: 04 अंक: 05 | मई -
: 10.55041/03470

एक अंतरराष्ट्रीय विद्वतापूर्ण || बहु-वर्षीय || ओपन एक्सेस || सभी प्रमुख डेटाबेस और मेटाडेटा में
अनुक्रमण

2025, (सभी अधिकार सुरक्षित) |.. |पेज
भाषा के लिए
समझ
बेहतर एनएलपी
मानक

इन्टरप्रेतगि
ध्यान
में
ट्रांसफार्मर
मॉडल
जेनकेल
एट
अल।

वविश
ध्यान
को
बढ़ाना
संरेखण
प्राप्त
++-स्तर
संरेखण
बना
बाहरी
संरेखण आंकड़ा
सामान्यीकृत नहीं
व्यापक एनएलपी को
कार्य

व्याख्या
शब्द
के लिए संरेखण
ट्रान्सफॉर्मर
डोऊ
और
न्यूबगि

हेड पूरनगि
और
देखरेख
संरेखण
सगिन्ल
का उत्पादन
व्याख्या
और
शुद्ध

शब्द
संरक्षण
आवश्यक है
अतिरिक्त
पर्यवेक्षण

ग्रेड-एसएम:
तस्वीर
स्पष्टीकरण
के लिए
ट्रांसफार्मर
फैसले
बर्कन एट
अल।

ग्रेडिएंट-
भारति
आत्मनर्णय
एमएपीएस
ऑफर
उन्नत
विविधता
के लिए
नमूना
फैसले
ध्यान केंद्रित
प्रमुख रूप से
पर
स्पष्टता

कैसे
करना
लल्म्स
अनुवाद करना?
इन्टरप्रेटिंग
साथ
पथ
पैच
झांग एट अल।
पाथ पैचिंग
ध्यान के लिए
प्रवाह अनुसंधान
केवल ~ 5%
ध्यान
प्रमुख कुंजी है
अनुवाद में
कार्य

अनुवाद में
ध्यान केंद्रित,
का अभाव
सामान्य
आंदोलन
अंतरदृष्टि

। अनुसंधान अंतराल

ध्यान-आधारित मॉडल और बहुभाषी दस्तावेज़ प्रसंस्करण में महत्वपूर्ण प्रगतियों के बावजूद, कई महत्वपूर्ण अनुसंधान अंतराल बने हुए हैं। सबसे पहले, जबकि ट्रांसफॉर्मर-आधारित मॉडल जैसे कि बिर्ट और जीपीटी उच्च गुणवत्ता की प्रदर्शक करते हैं

अनुवाद, वे अक्सर मूल लेआउट को संरक्षित करने और जटिल पीडीएफ दस्तावेजों के स्वरूपण की क्षमता में कमी करते हैं।

यह सीमा कानूनी, शैक्षणिक या वित्तीय दस्तावेजों जैसे संदर्भों में महत्वपूर्ण है, जहां संरचना बताती है अर्थ। दूसरा, अधिकांश अनुवाद प्रणाली क्लाउड-आधारित एपीआई पर बहुत अधिक भरोसा करती हैं, गोपनीयता की चिंताओं को पूरा करती हैं और सीमांकीकृत करती हैं

ऑफलाइन या सुरक्षित वातावरण में प्रयोज्यता। तीसरा, हालांकि विभिन्न अध्ययनों ने मॉडल को बढ़ाया है ध्यान वजुअलाइजेशन और ग्रेडिएंट विश्लेषण का उपयोग करते हुए, अभी भी सीमांकीकृत है कि कैसे विशिष्ट ध्यान प्रमुख विभिन्न वाक्य रचना संरचनाओं के साथ भाषाओं में अनुवादों को प्रभावित करते हैं। इसके अतिरिक्त,

मौजूदा उपकरण जैसे कम-रजिस्ट्रेशन या मल्टी-कॉलम पीडीएफ के साथ संघर्ष करते हैं,

जो डाउनस्ट्रीम को प्रभावित करता है

अनुवाद सटीकता। एक अन्य प्रमुख अंतर एंड-टू-एंड ओपन-सोर्स सिस्टम की अनुपस्थिति है जो ओसीआर, एनएलपी को एकीकृत करता है,

एक एकीकृत, मॉड्यूलर पाइपलाइन में अनुवाद, और लेआउट पुनर्निर्माण। अंत में, कई मॉडलों को प्रशिक्षित किया गया है

मुख्य रूप से उच्च-संसाधन वाली भाषाओं पर, दोनों के संदर्भ में कम-संसाधन वाली भाषा जोड़े को कम करना अनुवाद गुणवत्ता और मॉडल मूल्यांकन। इन अंतरालों को संबोधित करने से अधिक मजबूत, व्याख्या करने योग्य का मार्ग प्रशस्त हो सकता है,

और समावेशी बहुभाषी दस्तावेज़ अनुवाद प्रणाली।

। प्रस्तावित कार्यप्रणाली

प्रस्तावित प्रणाली

बहुभाषी पीडीएफ दस्तावेजों को निकालने, अनुवाद करने और पुनर्निर्माण के लिए एआई-संचालित पाइपलाइन विकसित करना है

उनके मूल स्वरूपण और संरचना को संरक्षित करते हुए। कार्यप्रणाली ऑप्टिकल चरित्र मान्यता को एकीकृत करती है

(), नेचुरल लैंग्वेज प्रोसेसिंग (), और ट्रांसफॉर्मर-आधारित ट्रांसलेशन मॉडल एक मॉड्यूलर फ्रेमवर्क में।

प्रक्रिया को पांच प्रमुख चरणों में विभाजित किया गया है:

3.1। पीडीएफ पाठ निष्कर्षण

और का उपयोग करना

, सिस्टम पहले फाइलों से मशीन-पठनीय और छवि-आधारित पाठ दोनों का पता लगाता है

और निकालता है। के लिए

इंटरनेशनल साइंटिफिक जर्नल ऑफ इंजीनियरिंग एंड मैनेजमेंट () : 2583-

वॉल्यूम: 04 अंक: 05 | मई -

: 10.55041/03470

एक अंतरराष्ट्रीय वदिवतापूर्ण | | बहु-वर्षिक | | ओपन एक्सेस | | सभी प्रमुख डेटाबेस और मेटाडेटा में अनुक्रमण

2025, (सभी अधिकार सुरक्षित) |.. | पेज

स्कैन या छवि-आधारित दस्तावेज़, - को छवियों को संपादन योग्य पाठ में परिवर्तित करने के लिए नियोजित किया जाता है जबकि

पाठ पदों, स्तंभों और पैराग्राफ की पहचान करना।

3.2। प्रीप्रोसेसिंग और भाषा का पता लगाना

नकाले

मूलपाठ

है

साफ, तार्किक ब्लॉकों में खंडित (जैसे, शीर्षकों, पैराग्राफ, टेबल), और भाषा का पता लगाने के माध्यम से पारित किया गया

स्रोत भाषा की पहचान करने के लिए एल्गोरिदम (जैसे,)। यह कदम सटीक भाषा-वशिष्ट सुनिश्चित करता है

अनुवाद हैडलिंग।

3.3। मशीन अनुवाद

पूर्वप्रभरण किया हुआ

पाठ का अनुवाद ट्रांसफार्मर-आधारित मॉडल जैसे कि या अनुवाद का उपयोग

करके किया जाता है, जो इस पर निर्भर करता है

परिनियोजन संदर्भ। सिस्टम लेआउट अखंडता को बनाए रखने के लिए टेक्स्ट ब्लॉक-बाय-ब्लॉक अनुवाद को संभालता है। वकिसति

मॉडल डोमेन-वशिष्ट भाषा उपयोग और बढ़ाया प्रवाह के लिए ठीक-ठीक हो सकते हैं।

3.4। लेआउट संरक्षण और पुनर्निर्माण

का उपयोग करते हुए

मेटाडेटा मूल दस्तावेज़ से (निर्देशांक, फ़ॉन्ट आकार, स्थिति), अनुवादित पाठ को एक नए में फिर से इंजेक्ट किया गया है

पीडीएफ लेआउट। या जैसे पुस्तकालय मूल दस्तावेज़ की संरचना को दोहराने में मदद करते हैं, जिसमें शामिल है

छवियाँ, टेबल और स्वरूपण।

3.5। पोस्टप्रोसेसिंग और मूल्यांकन

आउटपुट

पीडीएफ एक गुणवत्ता आश्वासन चरण से गुजरता है जिसमें समीक्षित स्थिरता जांच, लेआउट सत्यापन और शामिल है

या उल्का स्कोर का उपयोग करके भाषा की गुणवत्ता का मूल्यांकन। वैकल्पिक मानव-इन-द-लूप सत्यापन हो सकता है

उच्च-दांव अनुप्रयोगों के लिए एकीकृत।

वी। नष्टिकर्ष

यह

परियोजना,
 हमने बहुभाषी पीडीएफ अनुवाद के लिए एक एआई-आधारित प्रणाली का प्रस्ताव दिया जो अत्याधुनिक
 प्रौद्योगिकियों को एकीकृत करता है
 ऑप्टिकल चरित्र मान्यता (), प्राकृतिक भाषा प्रसंस्करण (), और ट्रांसफार्मर-आधारित
 मॉडल।द्वारा
 -, , और और जैसे शक्तिशाली अनुवाद मॉडल जैसे
 टूल का लाभ उठाना
 एपीआई का अनुवाद करें, सिस्टम कुशलता से अर्क, अनुवाद करता है, और बनाए रखते समय जटिल दस्तावेजों
 का पुनर्निर्माण करता है
 मूल लेआउट और संरचना।
 प्रस्तावित कार्यप्रणाली बहुभाषी दस्तावेज़ प्रसंस्करण में महत्वपूर्ण चुनौतियों को संबोधित करती है, जिसमें
 पाठ नष्टकरण भी शामिल है
 मशीन-पठनीय और छवि-आधारित पीडीएफ दोनों से, कई भाषाओं में सटीक अनुवाद, और
 स्वरूपण, तालिकाओं और छवियों का संरक्षण।इसके अतिरिक्त, सिस्टम की मॉड्यूलर और स्केलेबल प्रकृति
 अनुमत देता है
 अन्य दस्तावेज़ प्रसंस्करण वर्कफ़्लो के साथ आसान एकीकरण के लिए और इसकी क्षमताओं को अलग
 करने की क्षमता को अलग-अलग करने की क्षमता
 भाषा और डोमेन।
 हालांकि, भविष्य में वृद्धि के लिए अभी भी कई क्षेत्र हैं, जैसे कि कम-रजिस्ट्रेशन के लिए ओसीआर सटीकता में
 सुधार करना
 दस्तावेज़, कम-संसाधन भाषाओं के लिए अनुवाद की गुणवत्ता का अनुकूलन, और मॉडल व्याख्याता को बढ़ाने।
 कुल मिलाकर, यह परियोजना दस्तावेज़ प्रसंस्करण में भाषा की बाधाओं को तोड़ने के लिए एक व्यापक समाधान
 प्रदान करती है,
 बेहतर पहुंच और वैश्विक सहयोग को बढ़ावा देना।
 ।प्रतिक्रिया दें संदर्भ
 [१]।वासवानी, ए।, शेज़ीर, एन।, परमार, एन।, उस्ज़कोरटि, जे।, जोन्स, एल।, गोमेज़, ए।, कैसर, ।, और
 पोलोसुखनि, आई। (2017)।
 "ध्यान आप सभी की जरूरत है।"न्यूरोप्सि 2017 की कार्यवाही।
 [२]।बहडानौ, डी।, चो, के।, और बेगियो, वार्ड। (2014)।"संयुक्त रूप से संरेखित करने के लिए और संयुक्त रूप से
 सीखने के द्वारा तंत्रिका मशीन अनुवाद
 अनुवाद। " 2015 की कार्यवाही।
 [३]।लुओग, एम। टी।, फाम, एच।, और मैनिगि, सी। डी। (2015)।"ध्यान-आधारित तंत्रिका मशीन के लिए प्रभावी
 दृष्टिकोण
 अनुवाद। " 2015 की कार्यवाही।

इंटरनेशनल साइंटिफिक जर्नल ऑफ़ इंजीनियरिंग एंड मैनेजमेंट () : 2583-

वॉल्यूम: 04 अंक: 05 | मई -

: 10.55041/03470

एक अंतरराष्ट्रीय वद्वितापूर्ण || बहु-वर्षिक || ओपन एक्सेस || सभी प्रमुख डेटाबेस और मेटाडेटा में
 अनुक्रमण

2025, (सभी अधिकार सुरक्षित) |.. | पेज

[४]।जेनकेल, एस।, गेहरगि, जे।, और औली, एम। (2019)।"ट्रांसफार्मर मॉडल में ध्यान आकर्षित करना।"की
 कार्यवाही

एसीएल 2019।

[५]।बारकन, एस।, डागन, आई।, और शवार्ट्ज, वी। (2022)।"ग्रेड-एसएएम: ट्रांसफार्मर नरिणों के लिए दृश्य
 स्पष्टीकरण।"

2022 की कार्यवाही।

[६]।झांग, एक्स।, जेग, एक्स।, और ली, एस। (2024)।"एलएलएम कैसे अनुवाद करते हैं? पथ पैचिंग के साथ व्याख्या करना।"कार्यवाही

2024 का।

[[]]।डौ, जेड।, और न्यूबर्ग, जी। (2021)।"ट्रांसफार्मर के लिए व्याख्या योग्य शब्द संरेखण।" 2021 की कार्यवाही।

[[]]।स्मिथ, आर। (2007)।" इंजन।" 2007 की कार्यवाही।

[९]।वू, वाई।, शूस्टर, एम।, चेन, जेड।, ले।":

की न्यूरल मशीन अनुवाद प्रणाली। " 2016 की कार्यवाही।

[१०]।डेवलनि, जे।, चांग, एम। डब्ल्यू।, ली, के।, और टाउटनोवा, के। (2018)।"बर्ट: भाषा के लिए द्वादश ट्रांसफार्मर

समझ। " 2019 की कार्यवाही।