

Article

Revealing trends in geophysics using metadata and text analysis

Timofey Eltsov ^{1,†}, Maxim Yutkin ², Tadeusz W. Patzek ³

¹ Ali I. Al-Naimi Petroleum Engineering Research Center, King Abdullah University of Science and Technology; timofey_eltsov@kaust.edu.sa

² Ali I. Al-Naimi Petroleum Engineering Research Center, King Abdullah University of Science and Technology; maxim.yutkin@kaust.edu.sa

³ Ali I. Al-Naimi Petroleum Engineering Research Center, King Abdullah University of Science and Technology; tadeusz.patzek@kaust.edu.sa

* Correspondence: timofey_eltsov@kaust.edu.sa; Tel.: +966128087182

† Current address: 4700 KAUST, Thuwal, 23955-6900, Saudi Arabia

Version February 25, 2020 submitted to Geosciences

Abstract: Professional language evolution reveals the development of geophysics: researchers enthusiastically describe new methods of survey, data processing techniques, and objects of their study. Geophysicists publish their cutting-edge research at international conferences proceedings to share their achievements with the world. Tracking changes in the language allows one to identify trends and the current state of the science. Here, we describe the text analysis of the last 38 Annual Conferences organized by the Society of Exploration Geophysicists, one of the biggest geophysical gatherings. We split 24,500 articles into words and phrases and analyze the change in their usage frequency over time. We find that in 2019 the phrase “neural network” is used more often than “field data.” The word “shale” has become less commonly used, but the term “unconventional” is growing in occurrence. An analysis of conference materials and metadata allows one to identify trends in a specific field of knowledge and predict the development in the near future.

Keywords: geophysics; web data analysis; data mining; data analysis; text mining; words analysis;

1. Introduction

The last four decades showed a tremendous change in geophysics. An increase in computing power and technological progress allowed geophysicists to solve more and more complicated tasks. At the same time, the field of application of geophysics is expanding; the market of geophysical services is changing. We assume that a change of geophysical tasks, applications, geography, and technology will inevitably lead to a shift in the professional language. If one can track changes in the frequency of terms used in recent years, one can shed light on the current state of the industry and possibly predict future changes. Authors apply language processing methods to analyze changes in the professional language in geophysics.

The biases of different origin complicate big data [1]. In machine learning, the difference between training data set and test data set can cause biases. Massive sample study can lead to bias associated with an error resulting from sampling or study design [2]. Supposedly, it is better to have a smaller and more representative data set rather than a much bigger but biased data. We want to understand what the modern geophysical language looks like and what the future of geophysics will be. In this paper, we analyze only scientific articles presented at the Society of Exploration Geophysicists (SEG) Annual Conference and Exhibition. The committee selects the papers for the conference each year; this is the initial filtering. Also, it is worth noting that presenting at such a meeting is a demonstration of

30 the technical capabilities of industrial companies and the scientific viability of academic institutions.
31 Each annual conference proceedings is a cross-section of the state of geophysics, and we use it for
32 analysis and predictions.

33 The SEG Annual Conference and Exhibition is one of the biggest gatherings of geophysicists in
34 the world. Abstracts of the SEG Annual Conferences are a representation of the state of geophysical
35 science, approximated mainly to the oil and gas industry. Articles in the electronic version for the 38
36 years are available for analysis [3]. The SEG conducted all their Annual conferences in the USA, and
37 the last one was in San Antonio, TX. For analysis, the authors selected the proceedings of the SEG
38 Annual Conference, as the most representative set, that reflects state-of-art-technologies in geophysics.
39 Each conference proceedings is a reflection of the state of the industry in a particular year since, at this
40 event, both academic institutions and the industry present their best achievements in the field.

41 Besides conference proceedings, one can use journal articles for data mining as the volume of
42 the data for one year is comparable to the SEG Annual Conference and Exhibition Proceedings. The
43 number of publications per year is smaller, but they consist of full-size papers. However, the release of
44 articles in journals is carried out periodically, e.g., monthly or quarterly; at the conference, this happens
45 once a year. The research materials are usually published in journals and reported at conferences; the
46 proceedings include many of the results from full-sized articles. Moreover, the number of research
47 teams presenting their work is several times larger in the case of analysis of conference materials
48 compared to the study of one particular journal. SEG Annual Conference proceedings represent a
49 collection of scientific research from a large number of scientific teams in one place for each of the 38
50 years. This approach allows one to conduct a unique study and trace the dynamics of changes in the
51 industry.

52 **2. Materials and Methods**

53 We use open-source Python libraries: to transform, filter and process the text, and get metadata:
54 TextBlob, NLTK (Natural Language Toolkit), argparse, Pandas, Scrapy, Requests-HTML, sqlite3, and
55 NumPy. For printing the graphs, we use Matplotlib, Plotly, PIL (Python Imaging Library), and others.

56 We used digital versions of the SEG Annual Conference proceedings that have been available
57 online for 38 years. Fig. 1 shows the workflow scheme. Abstracts from the 1980s consisted of 1 or
58 2 pages; in the 1990s, it increased to four pages per abstract. We digitized articles in PDF format
59 from the SEG digital library website, converted them into plain TXT format using “pdftotext” with
60 “nopgbrk” (ignore page breaks), “enc ASCII7” (sets ASCII7 encoding for the output) and “eol” (sets the
61 end-of-line convention) flags. The text damp was filtered to remove common words, misspellings, etc.
62 from a NLTK dictionary, “stopwords.” After the initial filtering, we tokenized the text by year and
63 obtained si-, bi-, and trigrams¹. Further, we counted the number of times each word or phrase was
64 repeated. Finally, the entire text for 38 years is a three-dimensional array of words and phrases with
65 the corresponding number of repetitions for each year. We then analyze the list of words and phrases
66 during observation time and display the results in a graphical form.

67 While digitalizing abstracts from the 1980s, recognition errors, merged words, and typos occur.
68 Therefore, we present metadata analysis for the entire period; however, the phrase count is done only
69 for the period from 1990 to 2019.

70 In total, we analyzed 24,500 papers consisting of more than 57 million words or more than 383
71 million symbols.

¹ “sigram” - a word, “bigram” - two-word phrase, “trigram” - three-word phrase

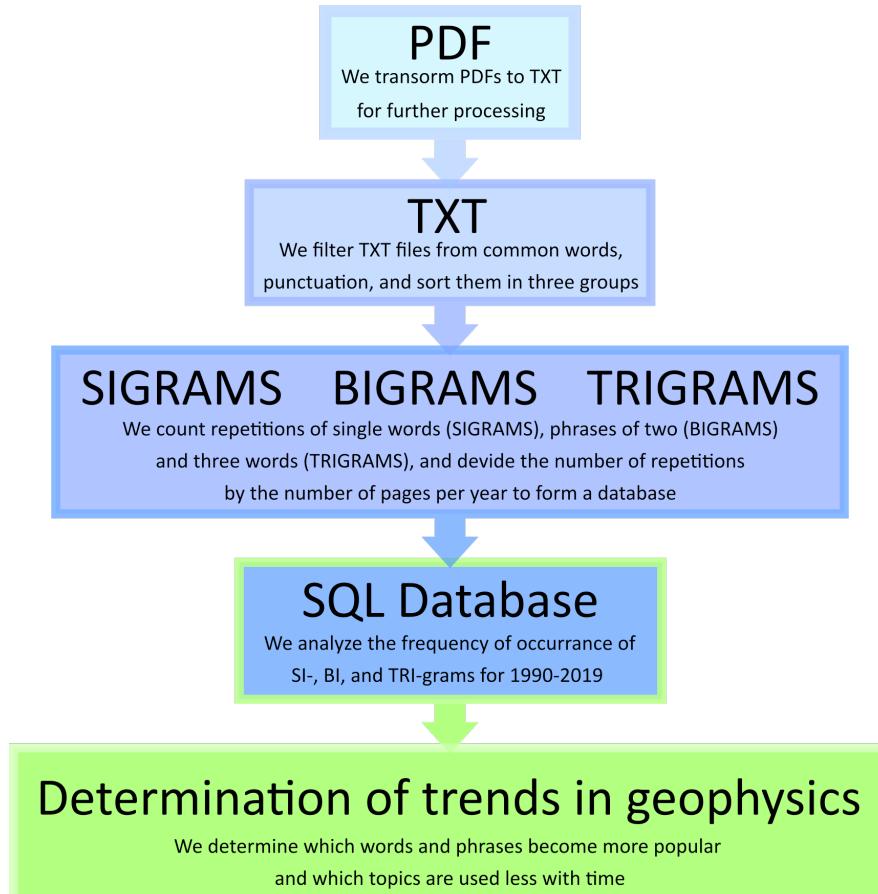


Figure 1. Data processing workflow.

72 3. Results

73 3.1. Technical terms analysis

74 In this section, we present our analysis of manuscript texts from the SEG Annual Conference.
 75 However, instead of focusing on average text length or sentence complexity, we look into the technical
 76 side. For example, we analyze and compare the frequency of occurrence of technical terms, such as
 77 "data" or "velocity." Such an analysis sheds some light on technology development and trends in the
 78 field.

79 The usage of the most frequently used English words ("the," "of," "and," "to," etc.) per page
 80 remains unchanged over the last decades. Therefore, we normalize the data to the number of pages of
 81 all articles every year. Often pages are not entirely filled with text; there are many graphs and formulas.
 82 Moreover, we know precisely the number of characters used, and we can estimate the number of pages.
 83 We calculate the average number of pages, Np , for each year using the formula: $Np_i = \frac{Ns_i}{3000}$, where i is
 84 the corresponding year, Ns - number of symbols, 3000 is the number of characters for the common one
 85 spaced web page. The estimated number of analyzed pages is 127.9 thousand. When analyzing the
 86 graphs in the present paper, one can state the number of times the phrase occurred per page each year.

87 3.1.1. Most common words and phrases

88 Fig. 2 shows the most commonly used words, two- and three-word phrases that appeared in
 89 conference materials from 1990 to 2019. The most frequent words are "data," "model," "velocity,"
 90 "seismic." Through the whole period of the study, the word "data" was mentioned more than 377700
 91 times, "seismic" 252400, "model" more than 251500 times, and "velocity" more than 223300 times in

92 thirty-eight years. For comparison, the mention of the word “that” was 324250 times. Most of the
 93 three- and two-word phrases are devoted to seismic exploration and seismic data processing.

94 Frequent use of these words tells us that most of the articles are about seismic exploration and
 95 seismic data processing. The terms “wellbore” and “logging” were more popular during the 1980s,
 96 and now their relative occurrence is declining.

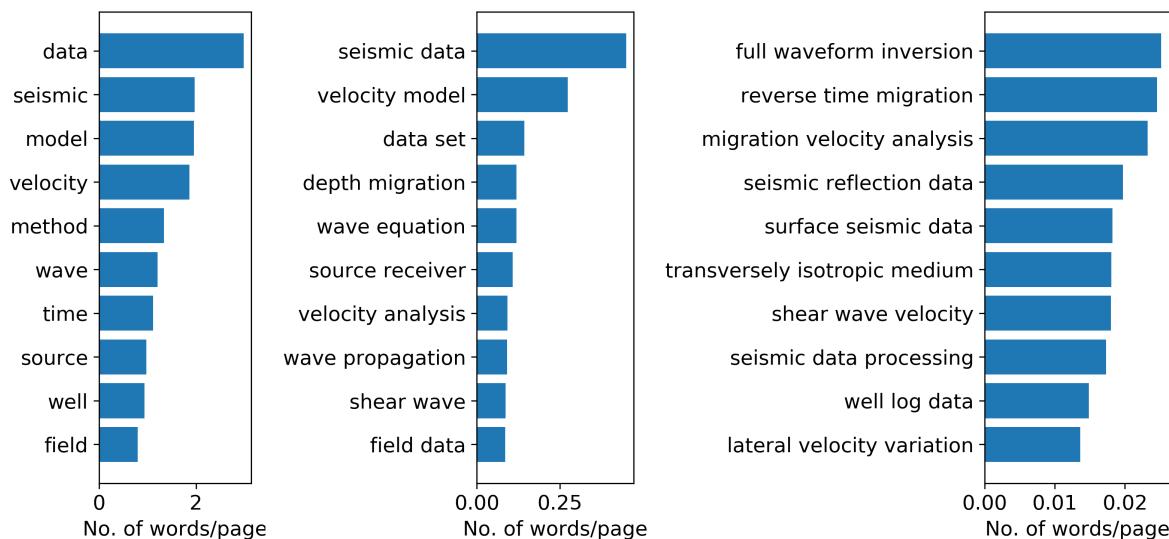


Figure 2. The average frequency of single words (left), two-word phrases (middle), and three-word phrases (right) per page for the most frequently used terms (the 1990s–2019). The total number of pages is 127.9 thousand.

97 While net average values provide information regarding the key concepts used over time, they
 98 are of less interest for absolutely the same reason. A more exciting approach is to monitor the evolution
 99 of other technical terms that constitute a subfield in geoscience or pertain to other disciplines. Such
 100 an analysis, however, is infinite. We limited the scope of this discussion to the objectives of the study,
 101 methods of data gathering and processing, shales, and neural networks. We also considered the fastest
 102 growing and declining trends in the publications.

103 3.1.2. Objects of the study

104 Fig. 3 breaks down the most used types of rocks. Each of the words on the left includes the
 105 most common names of rocks, e.g., sedimentary: shale, sandstone, conglomerate, carbonate, etc.;
 106 igneous: granite, diorite, basalt *etc.*; metamorphic: gneiss, phyllite, slate, *etc.* It shows the relative
 107 distribution of the objects of study: the majority of research deals with the sedimentary rocks. Terms
 108 that describe igneous rocks are used about ten times less than sedimentary, and the least used are
 109 metamorphic rocks related terms. The right part of Fig. 3 shows the occurrence of rock types with
 110 time. The shale revolution starting in 2007 is clearly notable. The most occurring names of rocks are
 111 “shale,” “sandstone” and “carbonate.” We note how “shale” peaks around 2015 and starts declining
 112 afterwards. Besides, there is a steady increase in the appearance of “carbonate” (the 1990s – 2005),
 113 while “sandstone” is used evenly over the years. An attentive reader will notice that during the growth
 114 of the use of the word “shale,” the fluctuation of the use of the words “sandstone,” and “carbonate”
 115 decreased.

116 We breakdown the sum of the names of geophysical methods used from 1990 to 2019, Fig. 4, left.
 117 They practically do not change over time, and we show the total in a pie chart. The whole pie chart is

the sum of all the words we use in Fig. 4². We show the occurrence of the most common geophysical methods, and it gives us an estimate of SEG Annual Conferences content. Three-fourths of the material relates to the collection and processing of seismic data; the remaining quarter accounts for all other methods. We see that the primary method discussed at the SEG Annual conference is “seismic,” its usage is an order of magnitude higher than other methods and it is still growing in the frequency of occurrence. It is worth noting, that the word “seismic” is mentioned about four times more often than the word “geophysics.” On the right part of Fig. 4, we breakdown the names of the main resources. We observe a slight increase in the frequency of the words “gas” and “water” from 1995 to 2015. Perhaps this is due to the increase in reservoir modeling related research. Figs. 3 and 4 show that for the last 30 years there are no significant changes in the use of geophysical methods and objects, with the exception of an increase in the frequency of occurrence of “shale” from 2010 to 2018 and a slight increase in the occurrence of “water” and “gas” from 1995 to present days. Significant changes occur in the use of terms related to specific methods of geophysical survey and data processing. This will be discussed in subsequent sections of the paper.

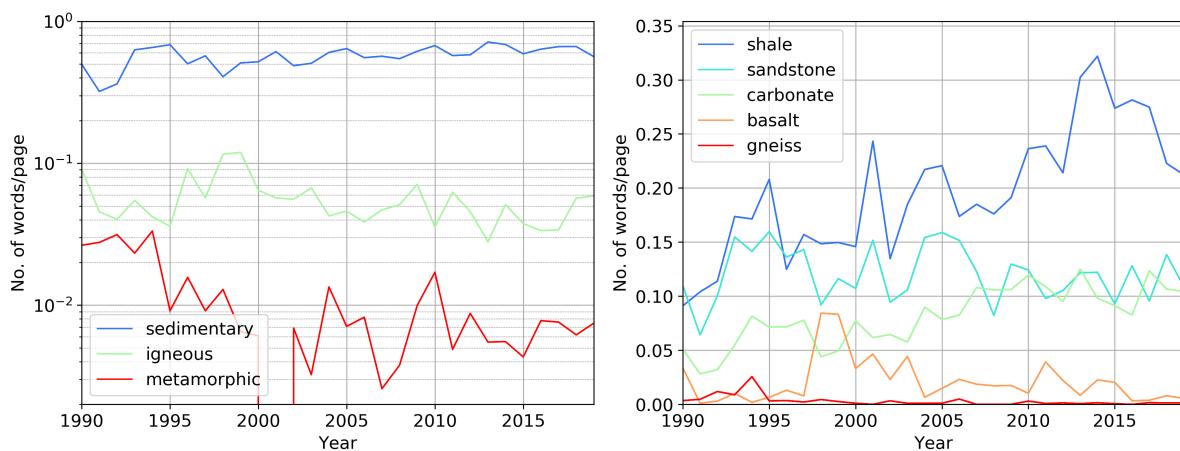


Figure 3. Frequency of use of different rock types (left) and most often used rock names (right). Rock types include most common rocks, e.g., sedimentary: shale, sandstone, carbonate; igneous: granite, diorite, basalt, etc.; metamorphic: gneiss, phyllite, slate, etc.

3.1.3. Processing and data acquisition methods

Of all the three-word phrases, the most frequently used now is “full waveform inversion” (Fig. 5), and it is still growing together with the abbreviation “FWI.” The second one is “reverse time migration,” the 2019 top three closes with “convolutional neural network.” Fig. 5 shows how the occurrence of “prestack depth migration” was substituted by “full waveform inversion” and “reverse time migration.” The frequency of occurrence is higher if we consider abbreviations. It is interesting to note that the abbreviations “FWI” and “RTM” are used more often than “PSDM” even when it was much more accessible. Perhaps this suggests a tendency to reduce and simplify the terms. The growth of the use of some terms inevitably supplants other words, provided that the volume of published material is approximately the same. While reviewing conference proceedings for the 38 years, we found many terms that were popular before, but do not find application in the modern world. Fig. 6, left, breaks down other trends in seismic data processing algorithms. We see that “machine learning” appears in the SEG Annual Conference proceedings more often in the past few years. The occurrence of the

² Each of the words represents the sum of the related words: “seismic,” “seismics”; “magnetic,” “geomagnetic,” “aeromagnetic”; “electromagnetic,” “em”; “gravity,” “gravimetry,” “gravimetric”; “electric,” “geoelectric”; “logging,” “borehole geophysics.”

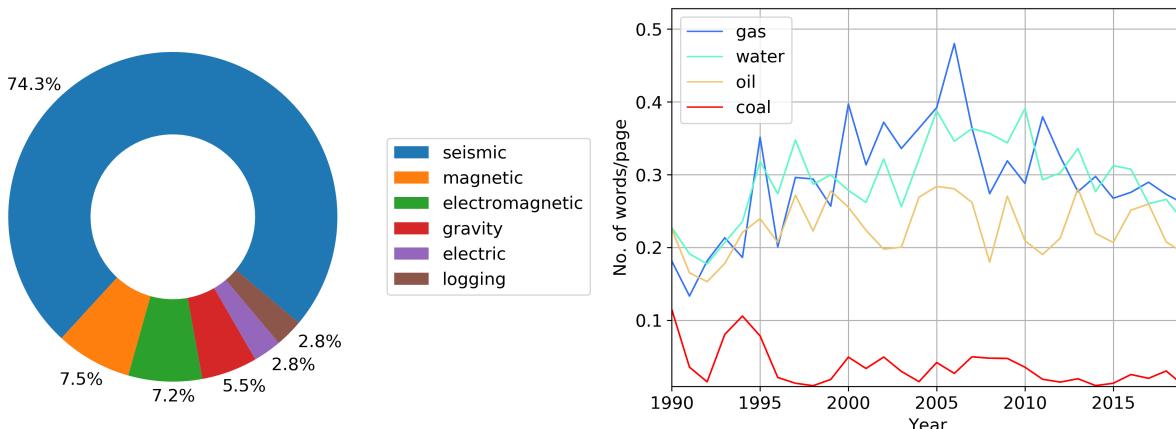


Figure 4. Geophysical methods of survey (left) and the most frequently used names of fossil fuels (right).

word “broadband” started to increase in the early 2010s with a corresponding decline in 2016-2018; it began to grow again in 2019. Using a wider frequency range and inclusion of the low frequencies proved to contribute to better resolution, penetration, and inversion [6]. Besides, we see an emergence in the rate of occurrence of the “Marchenko” method. “Marchenko” is a set of data-driven methods that help us to project surface seismic data to points in the subsurface, it relates Green’s function from a virtual source inside a medium to the reflection response at the surface of that medium [7,8]. “Markov”-chain-based approach is able to account for the change in seismic response of damaged structures [9], and it correlates with the occurrence of the word “seismicity.” The term “seismicity” is used for induced seismicity risk estimation [10], mine development [11], and other applications. It is known that “machine learning” and “neural networks” have recently significantly developed towards image recognition, and the in seismic data processing, and we will devote a separate part of the paper to this issue. Fig. 6, right, shows classic methods of seismic data processing and related terms. We see the rise and decrease in the appearance of these words in the last ten years, these methods were developed in the 1990s and now they have already been studied enough, and therefore their usage is decreasing. It should be noted that despite the decrease in the number of occurrence of the words “Kirchhoff” (migration), “CMP” (Common Mid Point) gather, “NMO” (Normal Moveout), “velocity analysis,” and “interferometry,” all of them are used in industrial seismic. These words are still used quite often, but the time of research and development of methods associated with these words was the 1990s and early 2000s. The decrease in the frequency of occurrence suggests that research on this topic has decreased.

3.1.4. Shale reserves

Fig. 7 shows the most often used names of shale plays on the left, and “fracking” (includes “hydraulic fracturing,” “frac,” and “fracking”) and “shale gas” + “gas shale” on the right. We observe peaking of shale-related terms from 2005 to 2015, which was followed by a decline in recent years. In the past 20 years, “Barnett” shale was mentioned more frequently than other shale deposits. In 2019 “Marcellus,” “Eagle” (Ford), and “Barnett” show the same occurrence, about one time per hundred pages. However, the term “fracturing” does not show such a fast decline. Despite the fact that the names of gas shale deposits reduced in use in the past three years, words that relate to the development and description of these deposits (“fracking,” “TOC” - total organic carbon, “unconventional”) do not show a decrease in use.

It is curious that in 2018, we observe an increase in the mention of the words “student,” “faculty,” and “researcher,” Fig. 8, left. Does this mean that the number of academic impacts has grown in the last year? You may notice the peaking of the word “engineer” after peaking of the word “student.” We

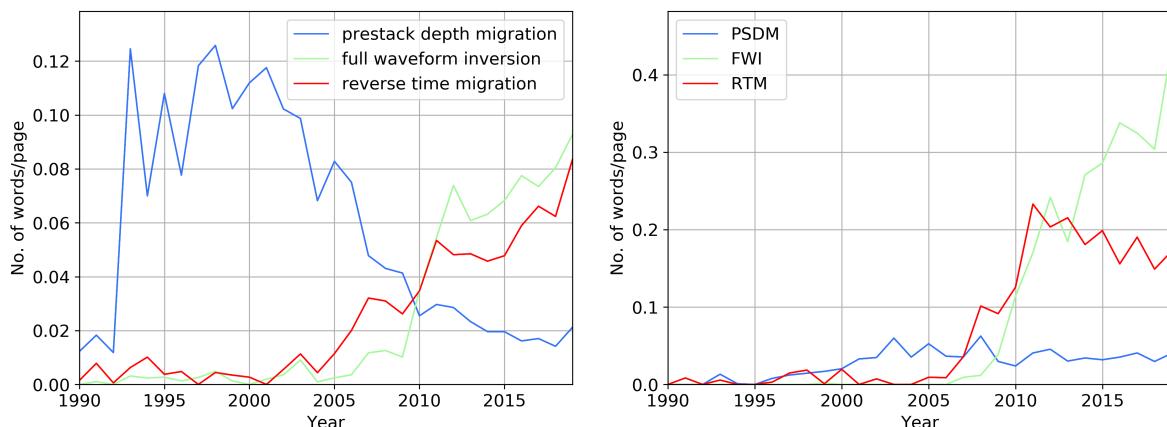


Figure 5. Change in the use of seismic data processing methods: full expression (left) and abbreviations (right).

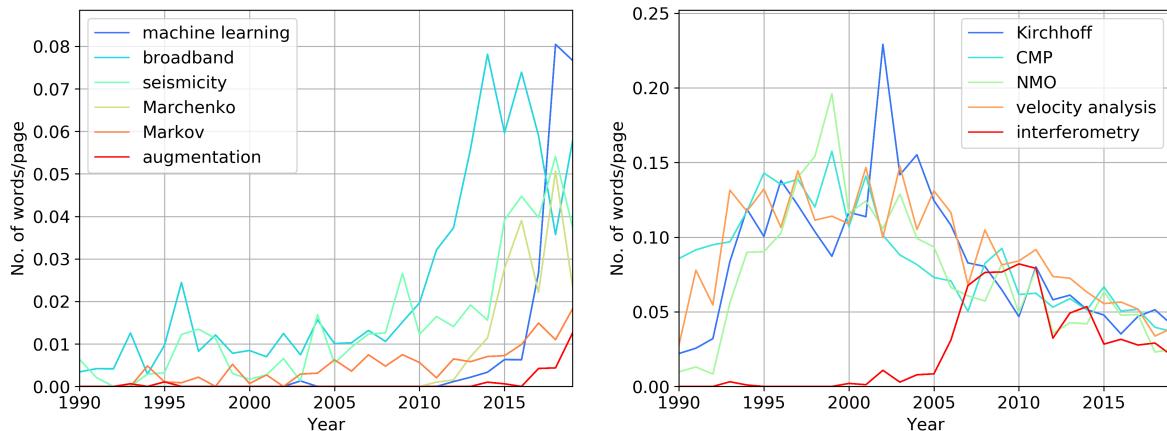


Figure 6. Trends in seismic data processing, terms, and algorithms that start to grow in usage (left) and decline in occurrence (right).

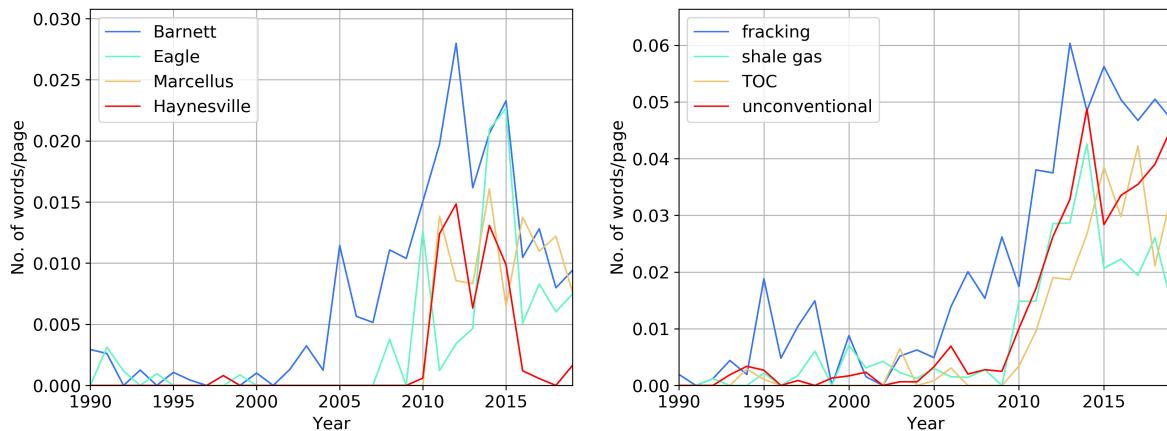


Figure 7. Most frequently mentioned shale reserves; change in usage of hydraulic fracturing and shale gas.

178 observe the growth in frequency of the word “researcher” word in the past ten years, and it appeared
 179 more often than engineer in 2019. During the 1990s, we see more of the word “engineer” in comparison
 180 to “researcher” and “scientist,” in the past decade, the situation has changed, bringing “researcher” to
 181 the first place. On the right side of Fig. 8, we observe an increase in the usage of the word “monitoring.”

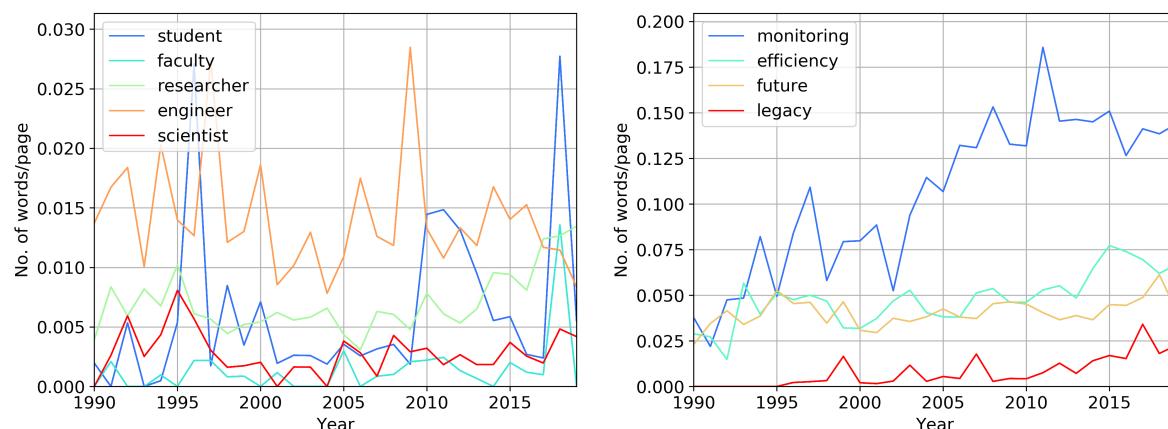


Figure 8. Change in words usage over time.

182 For example, this applies to microseismic monitoring and monitoring of the reservoir. The increased
 183 use of the term "monitoring" and "efficiency" indirectly indicates the concentration of researchers on
 184 the development of already explored deposits. The term "legacy" primarily refers to old data that is
 185 reprocessed using modern methods, including CNN. We used the word "future" regularly in the past
 186 30 years, perhaps, we can agree, the past is over.

187 3.1.5. Neural networks

188 We see that usually, the growth in the use of terms is saw-like; it is non-monotonic with individual
 189 peaks. Each peak represents the next phase of implementation, new research objects, and new teams
 190 that have mastered the method. "Neural networks" show a qualitatively different picture. From
 191 1990 to the beginning of 2000, attempts were made to use neural networks in geophysics, but they
 192 were suspended until 2016, in which a rapid growth in the use of this and related terms began. On
 193 average, we find a "neural network" phrase on every fourth page of the conference materials. If
 194 we observe an increased interest in this topic, then the researchers sincerely believe that using the
 195 methods of machine learning can solve many problems of geophysics. Given this context, we pose the
 196 question: Is the automation of geophysical data processing the main problem of modern geophysics?
 197 The authors believe that the main problem of geophysics is the lack of new research objects, such as
 198 hydrocarbon and other mineral deposits. Lack of survey objects is the reason for the increased interest
 199 in the development of methods for automatic processing of geophysical data. At the same time, the
 200 use of words "monitoring" and "efficiency" is growing, which indicates an understanding of the need
 201 for complete extraction of hydrocarbons and the monitoring of developed fields. Fig. 9 shows the
 202 appearance of "neural network," "deep learning," "artificial intelligence" and "field data," we use the
 203 last phrase for reference as it is always often used. In 2019, "neural network," occurred more often than
 204 "field data." It had already happened in 1993 and from 1999 to 2001, after that, it declined for a while,
 205 but now "neural network," "deep learning," and "artificial intelligence" have started to grow again
 206 ("artificial intelligence" appeared during the 1980s). The question is: Will the growth continue, or will
 207 it decline again like it did in 1993 -1995? The decline in interest in neural networks in early 2000 can
 208 be explained by an insufficient amount of computing power to realize the capabilities of the method.
 209 Now, technological progress allows us to use neural network methods successfully for face recognition.
 210 We also see attempts to introduce them to other areas of life. It is not necessary to be a rocket scientist
 211 to understand the reasons for the increasing interest in neural networks in geophysics. Experts want to
 212 automate geophysical data processing as much as possible. It remains only to understand whether
 213 we need to automate seismic data processing deeply. With time, we will have fewer oilfields to be
 214 explored, providing space for monitoring and increasing production efficiency.

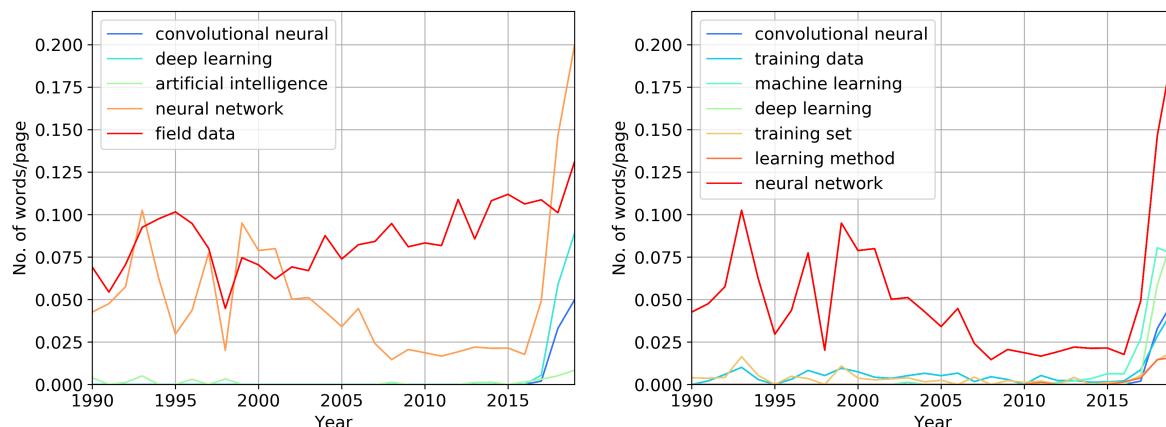


Figure 9. “Neural network” related two-word phrases. We display the phrase “field data” for the reference.

215 4. Discussion

216 The emergence of new techniques in geophysics inevitably leads to an increase in the use of terms
 217 from this field. The frequency of the occurrence of words can be used to track trends in the equipment,
 218 processing methods, math algorithms, and types of resources, including oilfields and the kind of rocks
 219 under study. The amount of hidden information is astounding. Such a study is unique because we
 220 have at our disposal a history of the development of geophysics. Moreover, it allows us to track exactly
 221 how the professional language changes over time.

222 It is interesting to know the terms that are gaining popularity now and discover the current trends
 223 in geophysics. Fig. 10 shows words with the highest growth in occurrence on the left and highest
 224 rate of decline on the right. As one can observe, the majority of words that have grown in occurrence
 225 relate to the neural network method. Is it possible to assume that these words will continue to gain
 226 popularity in the years ahead, and that the topic will remain relevant? For example, the phrases
 227 “streamer em” and “receiver deghosting” grew in occurrence at a very fast rate during 2011 – 2015,
 228 but since 2015, they decline as quickly as they were growing before. The word “fiber” and “fibre”³ is
 229 increasing in use almost as rapidly; this refers to fiber optics because seismic sensors based on fiber
 230 optics are now growing in use, showing their effectiveness in detecting faults filled with geothermal
 231 fluids [12], microseismic monitoring during hydraulic fracturing [13] and other applications. The term
 232 “distributed acoustic sensing” (DAS) shows good correspondence with the word “fiber” as a DAS
 233 is based on fiber-optics, and these terms are closely associated. Here, the use of the word is directly
 234 related to the production of the corresponding equipment. For “neural network,” one can use the
 235 existing computing power. Per contra, the development of optical fiber requires production. However,
 236 in 2019, we observe a decline in the usage of the word “fiber.” “Wasserstein” (metrics) and (data)
 237 “augmentation” have also grown in occurrence in the past three years but not that fast as “Marchenko.”
 238 Let us conclude, that the lack of research objects forces professionals to develop processing methods
 239 and, for example, reprocess legacy data. The picture on the right-hand side in Fig. 10 shows terms that
 240 decreased in occurrence in the past four years.

241 Interestingly, there has been a reduction in the use of the graphics processors by researchers as
 242 opposed to seven to eight years ago when the phrase was trending. “Barnett” shale is one of the
 243 most well studied, and the authors believe that the fading of interest in it is a natural phenomenon.
 244 Curiously, there was increased interest in basalt at the turn of the century, and we observe increased
 245 interest in the early 2010s.

³ “Chiefly British spelling of fiber” [15]

246 Besides “neural network” related terms (Fig. 11) on the left side, we observe an increase in usage
 247 of “tight sandstone” and “igneous rock.” It is interesting that for 30 years, “igneous rocks” were
 248 rarely discussed, except in 2009. In 2018 and 2019, we observe several papers discussing igneous rocks
 249 found on Chinese and Brazilian oil fields. Their acoustic and elastic properties must be considered
 250 in reservoir characterization [14]. On the right of Fig. 11 one can see two-word phrases that show a
 251 decrease in the frequency of occurrence in the past four years. When new research topics appear, new
 252 ones will partially or entirely replace old ones since the number of articles is limited every year.

253 Hill first described Gaussian beam migration in 1990. It is the seismic method that can image
 254 steeply dipping reflectors (more than 90 degrees) and will not produce unwanted reflections from
 255 the structure in the velocity model [16]. In 1993 at the SEG Annual Conference, we observe several
 256 papers reporting usage of beam migration in seismic data processing. In 2001 we notice an increase
 257 in the number of occurrence of “beam migration,” with the increase in computing capabilities, it
 258 became possible to use this method for 3D AVO analysis (Amplitude variation with offset) of small
 259 and medium-size 3D seismic surveys [17]. Interest in this method rises two more times, in 2008 and
 260 2015. Frequency peaks appear with enviable regularity every seven or eight years. Moreover, each
 261 subsequent peak is higher than the previous one. In 1990, a new method appeared; in 1993, we
 262 observe testing on synthetic data; in 2001, professionals report the results of processing small and
 263 medium volumes of data, in 2007 and 2008, the results of use on large objects in the Gulf of Mexico [18],
 264 CGGVeritas. For 25 years, we have seen the emergence of new technologies, testing, and application in
 265 field exploration. However, since 2015, we see a decrease in the rate of use of this term. Fig. 11 shows
 266 a reduction in the use of other seismic terms and Barnett shale.

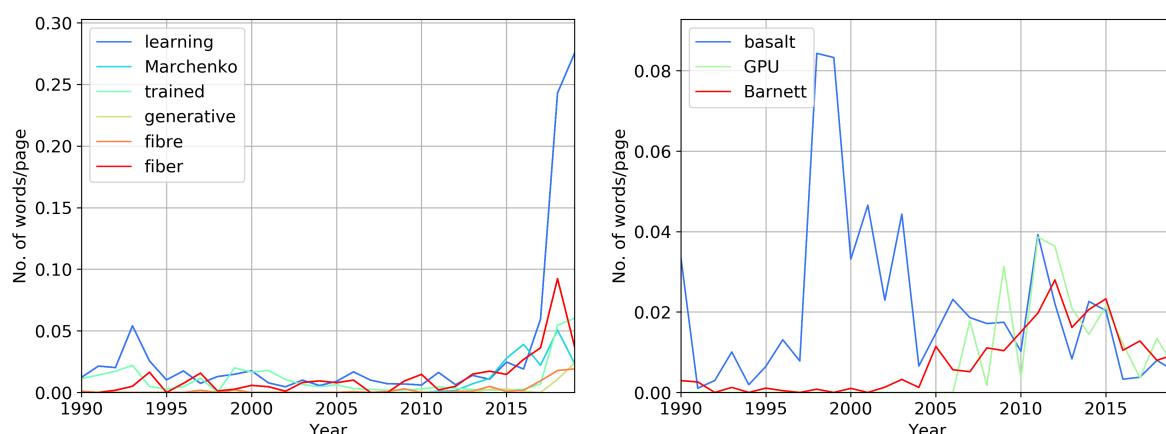


Figure 10. Words that show the highest rate of growth in occurrence (left) and decline (right) in the past four years.

267 Let us consider the most growing and declining three-word phrases, Fig. 12. “Convolutional
 268 neural network” (CNN) shows the fastest growth; the second one is “distributed acoustic sensing”
 269 (DAS), which is related to the fiber-optic measurement system. In the recent few years, researchers
 270 are using CNN to perform “seismic facies classification,” which is why we observe an increase in
 271 usage. We also see a relative increase for “ground penetration radar,” however, we see this term more
 272 often during the 1990s and early 2000. The right graph of Fig. 12 shows a decrease in the use of
 273 specific seismic terms as for the case of two-word phrases and the names of the shale deposits. From
 274 2010 to 2019, we observe an increase and decrease in interest in the phrase “towed streamer EM.”
 275 Towed streamer electromagnetic systems allow one to collect data at a high rate and over huge survey
 276 areas [19]. It is necessary to have significant objects to survey broad areas. Nowadays, there are less
 277 large-scale oil exploration projects, so the researchers use the corresponding terms less often.

278 It would be interesting to trace how the different methods are developing in geophysics, electrical
 279 exploration methods, petrophysics, engineering geophysics. For this reason, it is worthwhile to study

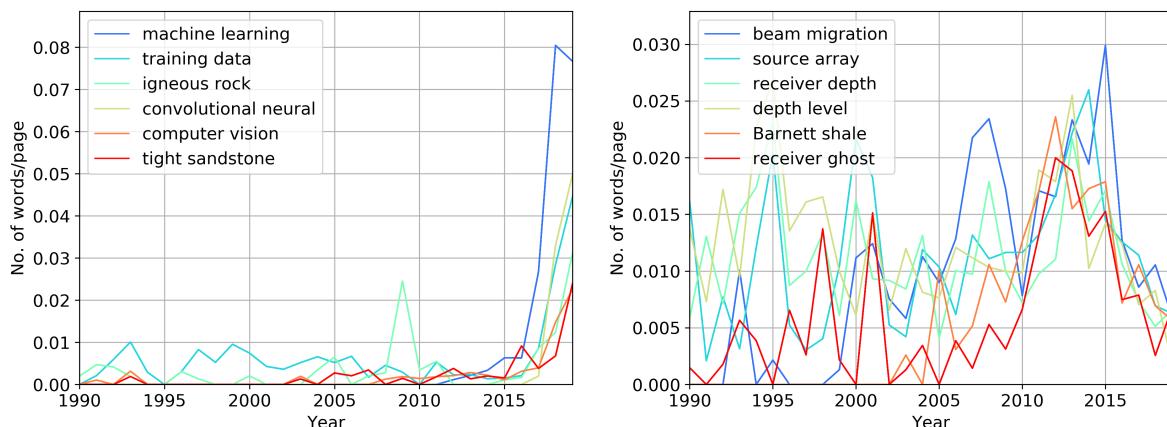


Figure 11. Two-word phrases that show the highest rate of growth in occurrence (left) and decline (right) in the past four years.

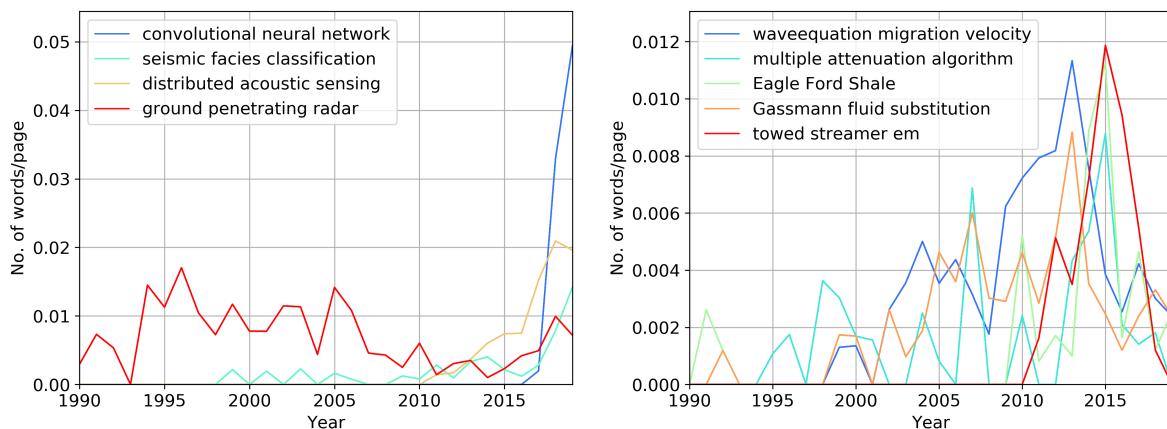


Figure 12. Three-word phrases that show the highest rate of growth in occurrence (left) and decline (right) in the past four years.

the materials of conferences and publications of other journals with a different specialization. Research of conference materials of other societies (SPWLA, EAGE, SPE) will provide a complete picture of the oil industry development.

The authors believe that the resulting database of industrial organizations and universities can be used to study the oilfield services and oil production market, or by students when searching for a place to study. We encourage readers to use our data available online [4]. The data includes the filtered word lists with the frequency of use each year, the number of pages, and the average number of co-authors. Thus, the reader will be able to conduct their research, test their hypotheses or assumptions.

5. Conclusions

We analyzed 24,500 papers, including 127,900 pages consisting of 57 million words, or more than 383 million symbols. Alteration in the professional language reflects the change in the industry and science. Over the 30 years, the objects and geophysical methods changed slightly. There has been an increased interest in "shales" in the last ten years. In the past six years, the frequency of the use of the word "shale" has been falling, but the use of the phrases "unconventional," "TOS," "hydraulic fracturing" has not decreased in recent years. At the same time, new methods of processing and capturing data appeared, and this led to a change in language. "Neural network" and related disciplines show the fastest growth in the last two years. The authors doubt that growth will continue at the same rate as the term "neural network" is already used more than "field data." More likely,

298 "neural network" related topics will occupy its niche in geophysics for the coming years. We see
299 an increase in the use of the words "Marchenko," "seismicity," and "broadband." We also observe
300 the rapid growth of the word "fiber," which is more likely related to fiber optic sensing systems.
301 Supposedly, we will see more projects on "monitoring" of oil and gas fields and increasing production
302 "efficiency," while there will be less work on the exploration of new oil and gas fields.

303 **Author Contributions:** Data mining and processing, software development, original draft preparation - Timofey
304 Eltsov; software development and analysis, review and editing of the draft - Maxim Yutkin; supervision, project
305 administration, historical analysis, review, and editing of the paper - Tadeusz W. Patzek.

306 **Funding:** Dr. Eltsov was supported by the KAUST Magnetic Sensor project, REP-2708.

307 **Acknowledgments:** Authors appreciate the responsiveness of the SEG team for permission to use digital data and
308 especially SEG Digital Publications Manager, Jeno Mavzer, for the useful advice and help. The authors are grateful
309 to their colleagues, and especially to Dr. Thomas Finkbeiner, for valuable and vital research recommendations.
310 The authors thank Dr. Sergey Yáskevich for consultations on exploration seismic. The authors are grateful to
311 Ilya Kolganov for the useful advice on the design of the graphs. We also would like to acknowledge Dr. Charles
312 Russell Severance for an informative Python course.

313 **Conflicts of Interest:** The authors declare no conflict of interest.

314 Abbreviations

315 The following abbreviations are used in this manuscript:

316 ASCII	American standard code for information interchange
AVO	Amplitude Variation with Offset
BP	BP plc., formerly The British Petroleum Company and BP Amoco
CGG	Compagnie Générale de Géophysique
CNN	Convolutional Neural Network
CMP	Common Mid Point
CSEM	The Controlled Source Electromagnetic
DAS	Distributed Acoustic Sensing
EAGE	European Association of Geoscientists and Engineers
EM	Electromagnetic
FWI	Full Waveform Inversion
GDP	Gross Domestic Product
GERD	Gross domestic Expenditure on Research and Development
GPU	Graphics Processing Unit
317 HTML	HyperText Markup Language
NLTK	Natural Language Toolkit
NMO	Normal Moveout
PDF	Portable Document Format
PIL	Python Imaging Library
PSDM	Prestack Depth Migration
RTM	Reverse Time Migration
R&D	Research and Development
SEG	Society of Exploration Geophysicists
SPE	Society of Petroleum Engineers
SPWLA	Society of Petrophysicists and Well Log Analysts
TXT	Text file
TOC	Total Organic Carbon
USA	The United States of America

318 References

- 319 1. Glauner, P.; Valtchev, P.; State, R. Impact of Biases in Big Data. In Proceedings of the European Symposium
320 on Artificial Neural Networks, Computational Intelligence and Machine Learning, Bruges, Belgium, 25-27
321 April 2018; pp 645–654.

- 322 2. Kaplan, R.; Chambers, D.A.; Glasgow, R.E. Big Data and Large Sample Size: A Cautionary Note on the
323 Potential for Bias. *CTS Journal* **2014**, *7*, 4, 342–346. doi:10.1111/cts.12178.
- 324 3. SEG Technical Program Expanded Abstracts. Available online: <https://library.seg.org/series/segeab>
325 (Accessed on 2 December 2019).
- 326 4. Eltsov, T. Data for SEG Annual Conferences analysis, 1982 - 2019. Available online:
327 https://github.com/ANPERC-source/SEG_Annual (Accessed on 2 February 2020).
- 328 5. Paper authorship goes hyper. Available online: <https://www.natureindex.com/news-blog/paper-authorship-goes-hyper> (Accessed on 17 October 2019).
- 330 6. Kroode, F.; Bergler, S.; Corsten, C.; Maag, J.W.D.; Strijbos, F.; Tijhof, H. Broadband seismic data — The
331 importance of low frequencies. *Geophysics* **2013**, *78*, 2, WA3–WA14. doi:10.1190/GEO2012-0294.1.
- 332 7. Lomas, A.; Curtis, A. An introduction to Marchenko methods for imaging. *Geophysics* **2019**, *84*, 2, 35–45.
333 doi:10.1190/geo2018-0068.1.
- 334 8. Thorbecke, J.; Slob, E.; Brackenhoff, J.; Neut, J.V.D.; Wapenaar, K. Implementation of the Marchenko method.
335 *Geophysics* **2017**, *82*, 6, WB29–WB45. doi:10.1190/geo2017-0108.1
- 336 9. Iervolino, I.; Giorgio, M.; Chioccarelli, E. Markovian modeling of seismic damage accumulation. *Earthquake
337 Engineering & Structural Dynamics* **2016**, *45*, November 2015, 441–461. doi:10.1002/eqe.
- 338 10. Weir, R.; Lines, L.; Lawton, D.; Eyre, T. The Duvernay Formation : the application of structure
339 and simultaneous inversion for reservoir characterization and induced seismicity. In Proceedings
340 of the SEG Annual Conference and Exhibition, Anaheim, USA, 14-19 October 2018; pp 2372–2376.
341 doi:10.1190/segam2018-2980345.1.
- 342 11. Barthwal, H.; Baan, M.V.D. Causative mechanism of microseismicity recorded in an underground mine.
343 In Proceedings of the SEG Annual Conference and Exhibition, Anaheim, USA, 14-19 October 2018; pp
344 2962–2966. doi:10.1190/segam2018-2980345.1.
- 345 12. Trainor-Guitton, W.; Jreij, S.; Guitton, A.; Simmons, J. Fault classification from 3D imaging of vertical DAS
346 profile. In Proceedings of the SEG Annual Conference and Exhibition, Anaheim, USA, 14-17 October 2018;
347 pp 4664–4668.
- 348 13. Chakraborty, G.; Chakraborty, D. Detecting microseismic events in downhole distributed acoustic sensing
349 data using convolutional neural networks. In Proceedings of the SEG Annual Conference and Exhibition,
350 San Antonio, USA, 15-20 September 2019; pp 4864–4868.
- 351 14. Penna, R.; Araújo, S.; Geisslinger, A.; Sansonowski R.; Oliveira, L.; Rosseto J.; Matos, M. Carbonate and
352 igneous rock characterization through reprocessing, FWI imaging, and elastic inversion of a legacy seismic
353 data set in Brazilian presalt province. *The Leading Edge* **2019**, *38*, 1, 11–19. doi:10.1190/tle38010011.1.
- 354 15. Merriam-Webster online dictionary. Available online: <https://www.merriam-webster.com/> (Accessed on 16
355 February 2020).
- 356 16. Hill, N.R.; Gaussian beam migration. *Geophysics* **1990**, *55*, 11, 1416–1428. doi:10.1190/1.1442788
- 357 17. Huang, S.; Sherrill, F.; Sengupta, M.K. Merits of amplitude preserving Kirchhoff beam migration method for
358 3D AVO analysis. In Proceedings of the SEG Annual Conference and Exhibition, San Antonio, USA, 9-14
359 September 2001; pp 1–4.
- 360 18. Ting, C.O.; Wang, D. Controlled beam migration applications in Gulf of Mexico. In Proceedings of the SEG
361 Annual Conference and Exhibition, Las Vegas, USA, 9-14 November 2008; pp 368–372.
- 362 19. Zhdanov, M.S.; Endo, M.; Sunwall, D.; Mattsson, J. Advanced 3D imaging of complex geoelectrical structures
363 using towed streamer EM data. In Proceedings of the SEG Annual Conference and Exhibition, New Orleans,
364 USA, 18-23 October 2015; pp 904–908.