

1 Article

# 2 A DCVS Reconstruction Algorithm for Mine Video 3 Monitoring Image Based on Block Classification

4 Xiaohu Zhao<sup>1,2</sup>, Xueru Shen<sup>1,2,\*</sup>, Kuan Wang<sup>1,2</sup> and Wanmei Li<sup>1,2</sup>

5 <sup>1</sup> The National and Local Joint Engineering Laboratory of Internet Application Technology on Mine, Xuzhou  
6 221000, China; TS16060112A3@cumt.edu.cn

7 <sup>2</sup> China University of Mining and Technology, Xuzhou 221000, China; TS16060247P3@cumt.edu.cn

8 \* Correspondence: TS16060247P3@cumt.edu.cn; Tel.: +86-151-525-01697

9

10 **Abstract:** Aiming at the problems that large amount of video monitoring image data in underground  
11 coal mines leads to difficulties in transmission and storage, compressed sensing theory is introduced  
12 to encode and decode video images, and a new distributed video coding scheme is proposed. In  
13 order to obtain more sparse representation and more general applicability, a block-based adaptive  
14 sparse base scheme is proposed. For the acquisition of side information, fixed weight is usually used  
15 to synthesize side information and the correlation between different image blocks is neglected, a  
16 block-based classification weighted side information generation scheme is proposed. Experimental  
17 results show that the block-based classification codec scheme can make full use of inter-frame  
18 correlation. Under the appropriate sampling rate, the PSNR value of video reconstruction increases,  
19 which effectively improves the quality of video frame reconstruction.

20 **Keywords:** compressed sensing; distributed video codec; sparse representation; side information  
21 reconstruction

22

---

## 23 1. Introduction

24 Coal is an important basic energy source in China and plays a decisive role in the development  
25 of economic construction [1]. As a non-renewable resource, the shallow coal mine resources are  
26 gradually mined and turned to the exploitation of deep coal seams, and the dangers that come with  
27 it are also increasing. Mine smart monitoring, coal seam identification and mineral detection have  
28 high requirements on the quality and real-time performance of video images, especially in  
29 unattended working areas in coal mines [2]. The amount of data of the monitoring image is large, so  
30 the requirement for the device in terms of image transmission and storage is high. The wireless  
31 sensor network (WSN) relies on low-power wireless sensor nodes, overcomes the difficulties of  
32 wired transmission wiring, incomplete detection information, high deployment costs and poor  
33 flexibility, and is widely used in the transmission of information in coal mines. Due to the  
34 particularity of underground coal mine working environment, a large number of sensor nodes need  
35 to be arranged downhole to monitor various signals such as gas concentration, temperature,  
36 humidity, and mine quake signals. Therefore, the amount of data the sensor needs to process is  
37 huge. However, the node energy and transmission bandwidth of the sensor are limited, and the  
38 energy supply of the mine node is difficult, so it is hard to apply a large amount of data collection.  
39 How to reduce the amount of information transmission and energy consumption of nodes is an  
40 urgent task to extend the life cycle of wireless network nodes, and data compression is a data  
41 processing technology that can effectively reduce the amount of data transmission. Compared with  
42 the traditional Nyquist sampling mode, Compressed Sensing [3-4] technology is devoted to  
43 searching for sparse solutions of underdetermined linear systems. The signal can be reconstructed

44 with a sampling rate much lower than the Nyquist sampling rate [5], which is suitable for signal  
45 processing.

46 In the video codec scheme, traditional video coding such as MPEG/H.26X, there will be a large  
47 amount of motion estimation and motion compensation [6], which makes the encoding side much  
48 more complex than the decoding side, this asymmetric coding scheme is suitable for a situation  
49 where multiple encodings are encoded at one time, such as broadcasting. However, more scenarios  
50 require low complexity, low power consumption, and low compression ratio at the encoder end. The  
51 distributed video coding (DVC) has low coding complexity and is suitable for use in scenarios  
52 requiring low complexity. In a distributed video coding scheme, the decoder performs motion  
53 estimation and compensation to find the correlation between adjacent frames, which makes the  
54 coding side less complex than the decoding side. When using CS theory to reconstruct the signal, the  
55 encoding side is less complex than the decoding side. However, using only the CS theory cannot  
56 provide us with a sufficiently low sampling rate, so we combine the CS theory and the DVC theory  
57 to form a new theory -- distributed compressed video sensing (DCVS) to make up for this deficiency.  
58 The DCVS technology allows independent encoding of multiple statistically related signals at the  
59 encoder and joint decoding at the decoder, which eliminates the need for complex motion estimation  
60 and compensation at the encoder end, thereby reducing the complexity of the encoder.

61 Combining the compressed sensing theory and the characteristics of distributed video codec,  
62 this paper proposes a block-based adaptive sparse representation and weighted side information  
63 reconstruction scheme. Firstly, the related theories of compressed sensing and distributed video  
64 coding are briefly introduced. Based on the correlation between video frames, a block-based  
65 classification sparse representation and classification weighted side information reconstruction  
66 scheme is proposed. Experimentally, the reconstructed PSNR value and time complexity are  
67 obtained, the experimental results are analyzed to verify the effectiveness of the algorithm.

## 68 2. Research Background Theoretical Analysis

### 69 2.1. Analysis of Characteristics of Wireless Sensor Network and Video Image in Underground Coal Mine

70 Under the special working environment of coal mines, if we want to fundamentally solve the  
71 problem of safe and efficient production of coal mines, the coal mining industry must shift from  
72 labor-intensive to technology-intensive, making it a new industry, new business model, and new  
73 model with high-tech features, and take the road of smart, few people (unmanned) and safe mining  
74 [7]. Under the transition of technology-intensive, wireless sensor networks are widely used in  
75 information transmission in coal mines because of their advantages such as low-power sensor nodes.  
76 The communication distance between WSN nodes is limited, and a large number of sensor nodes  
77 need to be arranged. However, the energy of the nodes is limited. The more data the nodes send and  
78 receive, the greater the energy consumption. But the downhole power supply in coal mines is  
79 difficult, so the data transmission needs to be reduced to reduce energy consumption. Large-scale  
80 deployment of sensor nodes will cause different nodes in the network to transmit similar data,  
81 resulting in a large amount of data redundancy and reducing transmission efficiency. Therefore, it is  
82 very necessary to apply compressed sensing theory to deal with video images in coal mines. In the  
83 collection of video images in underground coal mines, the collected images are mostly grayscale  
84 images composed of black and white. The hue is relatively single, but the correlation between  
85 adjacent images is very strong, and there is a large amount of redundant information between video  
86 frames. Compressed sensing technology can compress image information better and reduce  
87 transmission pressure.

### 88 2.2. Compressed Sensing Theory

89 The compressed sensing theory proposed by Donoho et al. is a novel sampling theory. From the  
90 perspective of analog to digital conversion, the compressed sensing theory provides a method of  
91 sampling at a rate lower than the Nyquist sampling rate without distortion recovery [8]. Different

92 from the traditional sampling theorem, the theory points out that the precondition for the signal  
 93 sparse representation is that the signal is sparse. The signal  $x$  is thus converted to other domain  
 94 space, represented by the sparse signal  $y$ . Consider an image signal  $x$  of length  $N$  whose  
 95 transform coefficient on some orthogonal base  $\Psi$  is expressed as:

$$96 \quad x = \Psi \theta \quad (1)$$

97 Where  $\theta$  is the coefficient vector of  $x$  on the orthogonal basis  $\Psi$ , the nonzero number of  $\theta$  is  $K$   
 98 and  $K \ll N$ , The signal  $x$  is called  $K$ -sparse signal. The linear prediction of the signal  $x$  is  
 99 performed using an observation matrix  $\Phi_{M \times N}$  ( $M \ll N$ ) that is incoherent with the orthogonal basis

100  $\Psi$ :

$$101 \quad y = \Phi x = \Phi \Psi \theta \quad (2)$$

102 Where  $y$  is the observation value and  $\Phi$  is the observation matrix. In this way,  $M$  linear  
 103 predictions which are much smaller than  $N$  are obtained, and by solving an optimization problem,  
 104 the original signals  $x$  can be reconstructed with high probability from these few projections. It can  
 105 be proved that such a projection contains enough information to reconstruct the signal. In this  
 106 theoretical framework, the sampling rate is not determined by the bandwidth of the signal, but  
 107 depends on the structure and content of the information in the signal.

### 108 *2.3 Distributed Video Coding Theory*

109 Distributed video coding (DVC) is a specific application of distributed source coding (DSC) in  
 110 video processing. DSC is the codec problem of distributed sources that are related to each other in  
 111 time or space. The information theory basis of DVC is Slepian-Wolf and Wyner-Ziv theorem [9].  
 112 Slepian-Wolf theory describes the conditions that need to be satisfied for lossless coding of related  
 113 sources. It is proposed that the related source "independent coding, joint decoding" can be the same  
 114 compression efficiency as "joint coding, and joint decoding". The Wyner-Ziv theory is supplement  
 115 and development based on the Slepian-Wolf theory and discusses the description of the  
 116 rate-distortion function  $R_{X/Y}^{WZ}(D)$  of the relevant source coding in the case of lossy coding. The  
 117 Wyner-Ziv rate-distortion function shows the lower limit of the code rate for distributed coding  
 118 under a distortion constraint. In the case of lossy coding, the rate distortion function obtained by  
 119 using the side information only at the decoding end and using at the codec side is consistent [6].  
 120 Generally, Wyner-Ziv encoding can be equivalent to quantifying sources and then performing  
 121 Slepian-Wolf encoding.

## 122 **3. Video Image Sparsity Analysis and Compression Coding Processing**

### 123 *3.1. Sparse Representation of the Signal*

#### 124 *3.1.1. Block-prediction and DCT Mixture Sparse Basis*

125 From Fourier transform to wavelet transform to later multi-scale geometric analysis, the  
 126 purpose of the scientists' research is to study how to provide a more concise and direct analysis  
 127 method for signals in different functional spaces. All of these transformations are designed to exploit  
 128 the characteristics of the signal and sparsely represent it, or to improve the nonlinear function  
 129 approximation of the signal, and further studying the degree of sparsity of the signal or the degree of

130 energy concentration of the decomposition coefficients using a set of bases in a certain space. From  
 131 the perspective of signal analysis, the Fourier transform is the basis for signal and digital image  
 132 processing. Wavelet analysis brings signal and digital image processing to a whole new field.  
 133 Multi-scale geometric analysis is a new generation of signal analysis tool following wavelet analysis.  
 134 It has many excellent features such as multi-resolution, localization and multi-directionality. It is  
 135 more suitable for processing multi-dimensional signals such as images. All these studies laid the  
 136 foundation for the theory of compressed sensing.

137 In order to get a better sparse representation, many methods have been studied. In the  
 138 traditional DCVS framework, linear or orthogonal transform bases such as discrete cosine transform  
 139 (DCT) or discrete wavelet transform (DWT) are widely used due to their simple and efficient  
 140 features. Traditional video compression or distributed video coding (DVC) technology based on  
 141 motion estimation and motion compensation all rely on a high-cost mechanism, that is, perception  
 142 or sampling and compression are not performed continuously, which leads to waste of resources. In  
 143 other words, most of the collected raw video data may be discarded in the complex compression  
 144 process [10]. However, learning-based dictionaries can provide a more sparse representation of  
 145 image signals than these pre-specified base classes. A dictionary learned from a set of blocks globally  
 146 extracted from the previous reconstructed neighboring frames is used as sparse basis in [11].  
 147 Nevertheless, this dictionary is based on K-SVD algorithm [12] and can be only applied in some  
 148 specific applications. Dictionaries reported in [13] and [14] are generated by a linear combination of  
 149 neighboring blocks in preceding and following key frames. However, with limited reference  
 150 information, the learned dictionary may not be sufficiently redundant. In [15], a simple side  
 151 information generation technique based on correlation analysis of CS measurements is effective in  
 152 complexity reduction. But the rate-distortion performance is unsatisfying compared with other  
 153 schemes.

154 In the traditional video codec, a hybrid coding method combining prediction and DCT can  
 155 quickly reduce spatiotemporal correlation and obtain a more sparse representation of the video  
 156 frame. In order to obtain better rate-distortion performance, we propose a new block prediction and  
 157 hybrid sparse basis for non-key frames with reference to this hybrid coding method. As shown in  
 158 equation (3), the new sparse basis consists of two parts, the initial DCT matrix and the  
 159 block-prediction basis. In this way, initial DCT basis is acquired by linear transform in DCT domain  
 160 while block-prediction basis is based on the SI generated from the adjacent decoded key frames.

$$161 \quad [\Psi_{DCT}; \Psi_{inter}] = \Psi_i \quad (3)$$

162 Where  $\Psi_{DCT} \in R^{N \times N}$  is an initial DCT basis,  $\Psi_{inter} \in R^{N \times 1}$  is a block-prediction basis and

163  $\Psi_i \in R^{N \times (N+1)}$  is a newly built hybrid DCT basis.

### 164 3.1.2. Classification Judgment Standard

165 In order to obtain more sparse representations and more general applicability, this paper  
 166 proposes a new hybrid sparse scheme that combines linear DCT bases and block prediction bases to  
 167 solve sparse representation problems. Adaptive block-based prediction is used to generate side  
 168 information and to learn block prediction basis. Since the video is composed of consecutive frames

169 and the time redundancy is particularly large, that is, the inter-frame correlation is very strong. By  
 170 utilizing the correlation between adjacent frames, the scheme can achieve more sparse  
 171 representation while reducing the complexity.

172 Different regions in the video sequence have different inter-frame correlations. We propose an  
 173 adaptive block prediction scheme for sparse representation. Perform the non-overlapping block  
 174 processing of the reconstruction results of the two adjacent key frames before and after the current  
 175 non-key frame. The difference between the  $i$ -th block  $x_{t+1}^i$  of the following frame and the  $i$ -th block

176  $x_{t-1}^i$  of the previous frame at the corresponding frame may reflect the correlation between the  
 177 inter-frames in the video sequence. Consider the residual energy values of the two corresponding  
 178 frame sub-blocks as the classification criteria, and define the residual energy value as shown in  
 179 equation (4):

$$180 \quad E(x_{t-1}^i, x_{t+1}^i) = \|x_{t-1}^i - x_{t+1}^i\|_2^2 \quad (4)$$

181 However, for video frames where video scenes change rapidly, the residual energy value of the  
 182 corresponding position may be large. Therefore, simply using the residual energy value as the  
 183 classification criterion will cause the threshold value to be too dependent on the sequence itself,  
 184 making the algorithm less versatile. Based on this, it is thought that the ratio of the residual energy to  
 185 the energy block of the previous key frame is used as the classification criterion, as defined by  
 186 equation (5):

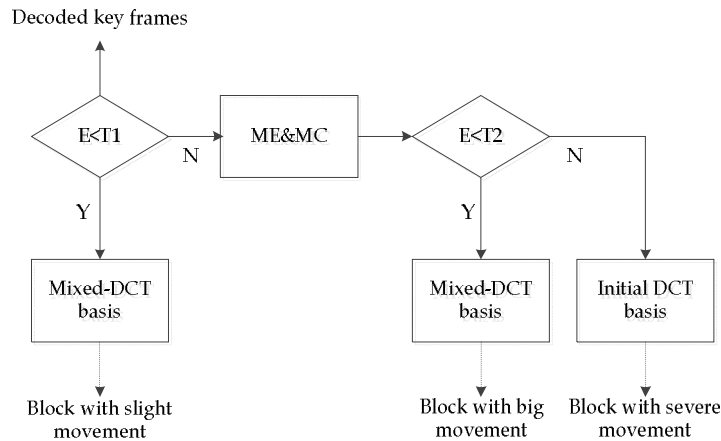
$$187 \quad e(x_{t-1}^i, x_{t+1}^i) = \frac{E(x_{t-1}^i, x_{t+1}^i)}{E(x_{t-1}^i)} \quad (5)$$

188 The video frame sub-blocks are classified using thresholds  $T_1$ ,  $T_2$  ( $T_1 < T_2$ ) for video sequences.  
 189 Calculate the value of the judgment function of the  $i$ -th block according to equation (5). If

190  $e(x_{t-1}^i, x_{t+1}^i) < T_1$ , then the block is defined as block with slight movement; if  $T_1 < e(x_{t-1}^i, x_{t+1}^i) < T_2$

191 , then the block is defined as block with big movement; and if  $e(x_{t-1}^i, x_{t+1}^i) > T_2$ , then the block is

192 defined as block with severe movement. In the experiment, the video frame is subjected to block  
 193 classification processing. The classification judgment thresholds  $T_1$  and  $T_2$  take 0.003 and 0.015  
 194 respectively. After block classification of non-key frames, the classification results and  
 195 corresponding sparse base strategies are shown in Figure 1.



196

197

Figure 1. Block classification of non-Key frames

198

199

200

201

202

203

204

205

206

207

208

209

After the non-key frames are blocked according to the defined classification criteria, different sparse-based strategies are proposed for different types of blocks. Algorithm 1 describes the process of adaptive hybrid DCT-based learning. For the block with slight movement, predictions are generated by generating the side information by performing forward and backward motion estimation on adjacent decoded key frames. The block prediction basis is learned and combined with the DCT basis to construct a hybrid DCT basis. For the block with big movement, by selecting different weights and combining forward and backward motion estimation to obtain side information, the construction method of the hybrid DCT base is similar to the block with slight movement. For the block with severe movement, the inter-frame correlation between adjacent frames is small, and motion estimation and compensation cannot accurately predict these fast-changing blocks. Therefore, for this type of block, only the initial DCT basis can be used. The description of algorithm 1 is as follows.

210

Algorithm 1: The hybrid DCT sparse basis algorithm based on block prediction.

---

**Input:**  $x_{t-1}^i, x_{t+1}^i$ : Neighboring decoded key frames as reference for the  $i$ -th non-key frames.

---

1. Calculate the residual energy value of the corresponding frame sub-block:

$$E(x_{t-1}^i, x_{t+1}^i) = \|x_{t-1}^i - x_{t+1}^i\|_2^2;$$

2. Calculate the ratio  $e(x_{t-1}^i, x_{t+1}^i) = \frac{E(x_{t-1}^i, x_{t+1}^i)}{E(x_{t-1}^i)}$  of the residual energy to the energy block

of the previous key frame as a classification criterion;

3. Calculate side information  $SI^i = \alpha x_{t-1}^i + (1 - \alpha) x_{t+1}^i$ , Where  $\alpha$  selects different weights according to the classification;

4. If  $e(x_{t-1}^i, x_{t+1}^i) < T_1$ ,  $\alpha = 0.8$ ,  $\Psi_{inter} = x_t^i / 255$ ,  $\Psi_t(i) = [\Psi_{DCT}; \Psi_{inter}]$ ;

5. If  $T_1 < e(x_{t-1}^i, x_{t+1}^i) < T_2$ ,  $\alpha = 0.5$ ,  $\Psi_{inter} = x_t^i / 255$ ,  $\Psi_t(i) = [\Psi_{DCT}; \Psi_{inter}]$ ; otherwise,
-

$$\alpha = 0.3, \quad \Psi_t(i) = \Psi_{DCT}.$$

**Output:**  $\Psi_t(i)$ : Hybrid DCT sparse basis for the  $t$ -th block in the  $i$ -th non-key frame.

### 211 3.2. Design of Observation Matrix

212 After the sparse representation, how to design a stable  $M \times N$  dimensional observation  
213 matrix  $\Phi$  that is not related to the transform basis  $\Psi$ , and ensure that the important information of  
214 the sparse vector  $\theta$  when descending from the  $N$  dimension to the  $M$  dimension is not destroyed,  
215 this is the second problem to be solved.

216 In compressed sensing theory, after the sparse coefficient vector  $\theta = \Psi^T X$  of the signal is  
217 obtained by transformation, the observation part of the compressed sample is designed. The  
218 purpose of the observation matrix design is how to sample and obtain  $M$  observations, and it can be  
219 ensured that a signal  $X$  of length  $N$  or a sparse coefficient vector  $\Psi$  of the base  $\Psi$  can be  
220 reconstructed from it. The importance of the observation matrix design is that if the signal  $X$  is  
221 destroyed during the observation process, it is impossible to reconstruct signal successfully. The  
222 observation process actually uses the  $M$  row vectors  $\{\varphi_j\}_{j=1}^M$  of the  $M \times N$  dimensional  
223 observation matrix  $\Phi$  to project the sparse coefficient vector. That is, to calculate the inner product  
224 between  $\theta$  and each observation vector  $\{\varphi_j\}_{j=1}^M$ , and get  $M$  observations  
225  $y_j = \langle \theta, \varphi_j \rangle (j = 1, 2, \dots, M)$ , for the observation vector  $Y = (y_1, y_2, \dots, y_M)$ , that is

$$226 \quad Y = \Phi \theta = \Phi \Psi^T X = A^{CS} X \quad (6)$$

227 The sampling process here is adaptive, that is,  $\theta$  does not have to change according to the  
228 change of  $X$ , and the observed data is no longer a point sampling, but a more general  $K$  linear  
229 functional of the signal.

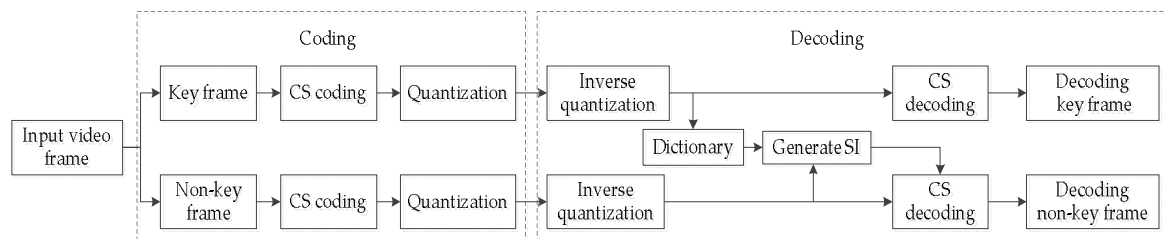
230 Given a vector  $Y$ , finding  $\theta$  from equation (6) is a linear programming problem. However,  
231 since  $M \ll N$ , that is, the number of equations is less than the number of unknowns, this is an  
232 underdetermined problem and the deterministic solution cannot be calculated. However, if  $\theta$  has  
233  $K$ -sparseness ( $K \ll M$ ), it is expected to find a certain solution just try to determine the appropriate  
234 position of the  $K$  non-zero coefficients  $\theta_i$  in  $\theta$ . Since the observation vector  $Y$  is a linear  
235 combination of  $K$  non-zero coefficient column vectors, a linear equation of  $M \times N$  is formed to  
236 solve the specific values of these non-zero terms. For whether there are necessary and sufficient  
237 conditions for determining the solution, Restricted Isometry Property (RIP)[16] gives the definition,  
238 The necessary and sufficient conditions for the definition are consistent with the geometric  
239 properties proposed by Candes and Tao et al that the sparse signals must maintain under the  
240 observation matrix. In the reconstruction process of the signal, in order to completely reconstruct  
241 the signal, it must be ensured that the observation matrix does not map two different  $K$ -term  
242 coefficient signals into the same sampling set, which requires that the matrix formed by the  $M$   
243 column vectors extracted from the observation matrix is non-singular. So the key to solving such  
244 problems is how to determine the position of the non-zero coefficients to form a linear system of  
245 equations.

246 How to judge the RIP property of a given matrix  $A^{CS}$  is a combination complexity problem.  
247 Finding an alternative method that can easily implement RIP characteristics is the key to  
248 successfully constructing the observation matrix. If the observation matrix  $\theta$  and the sparse basis  
249  $\Psi$  are guaranteed to be uncorrelated,  $A^{CS}$  can largely satisfy the RIP property. Irrelevant means  
250 that the vector  $\{\varphi_j\}$  cannot be sparsely represented by  $\{\Psi_j\}$ . The stronger the irrelevance, the  
251 more coefficients are needed to represent each other; on the contrary, the stronger the relevance.

252 Candes E J et al. proved that when  $\Phi$  is a Gaussian random matrix, the incoherence and RIP  
 253 conditions can be satisfied with a large probability. Therefore, the observation matrix  $\Phi$  in this  
 254 paper adopts a random Gauss matrix. A random Gaussian matrix has a useful property: for a  
 255  $M \times N$  random Gaussian matrix  $\Phi$ , it can be shown that when  $M \geq cK \log(N/K)$ ,  $\Phi\Psi^T = A^{CS}$   
 256 satisfies the RIP property with a large probability (where  $c$  is a small constant). Therefore, a  $K$ -term  
 257 sparse signal length of  $N$  can be reconstructed with high probability from  $M$  observations. The  
 258 random Gaussian matrix is not related to the matrix formed by most fixed orthogonal bases. This  
 259 property determines the choice that we choose it as an observation matrix. When other orthogonal  
 260 bases are used as the sparse transform base, the RIP can be satisfied.

### 261 3.3. Compressed Codec Processing Research

262 Compressed sensing is a new way of signal acquisition that allows us to design very simple  
 263 video encodings that can be implemented on mobile devices with limited resources. However, the  
 264 CS-based video codec proposed by the predecessors either requires a conventional video encoder or  
 265 a feedback channel, which increases the complexity of the codec. The distributed compressed video  
 266 sensing (DCVS) codec scheme proposed in this paper uses CS only at the encoding end, and the  
 267 decoder uses a correlation between CS measurements of adjacent frames to form a new side  
 268 information generating scheme. The new DCVS block diagram we proposed is shown in Figure 2.  
 269 This codec is completely CS-based and does not involve traditional video encoders. Both key and  
 270 non-key frames are encoded as CS measurements and no feedback channel is required.



271

272

Figure 2. DCVS video codec scheme

273 At the encoding end, the video sequence is divided into a number of Group of Pictures (GOPs).  
 274 Video frames are divided into key frames and non-key frames, non-key frames also called  
 275 Wyner-Ziv (WZ) frames, which form a group of pictures. Each GOP contains one key frame and  
 276 several non-key frames. Both key frames and non-key frames are coded with similar CS theory.

277 At the encoding end, each key frame is reconstructed by the SAMP algorithm, which redefines  
 278 the  $l_1$  minimization problem, as shown in equation (7).

$$279 \min_{\alpha_k} \frac{1}{2} \|y_k - A\alpha_k\|_2^2 + \tau \|\alpha_k\|_1 \quad (7)$$

280 Where  $y_k$  is the CS measurement of the key frame obtained at the decoding end,  $A = \Phi\Psi$  has

281 been described in the previous,  $\alpha_k \in R^N \times 1$  is the sparse coefficient vector of the solution. Key

282 frame  $\hat{x}_k$  is obtained by  $\hat{x}_k = \Psi\hat{\alpha}_k$ , where  $\hat{\alpha}_k$  is the optimal solution of  $\alpha_k$  in equation (7).

283 The decoding of non-key frames are assisted by the side information generated by the  
 284 dictionary and the side information is generated by the inverse quantized CS measurement of the  
 285 key frame. The effect of side information is only apparent when there is sufficient correlation  
 286 between the measured values of the key frame and the WZ frame. In a video sequence, adjacent

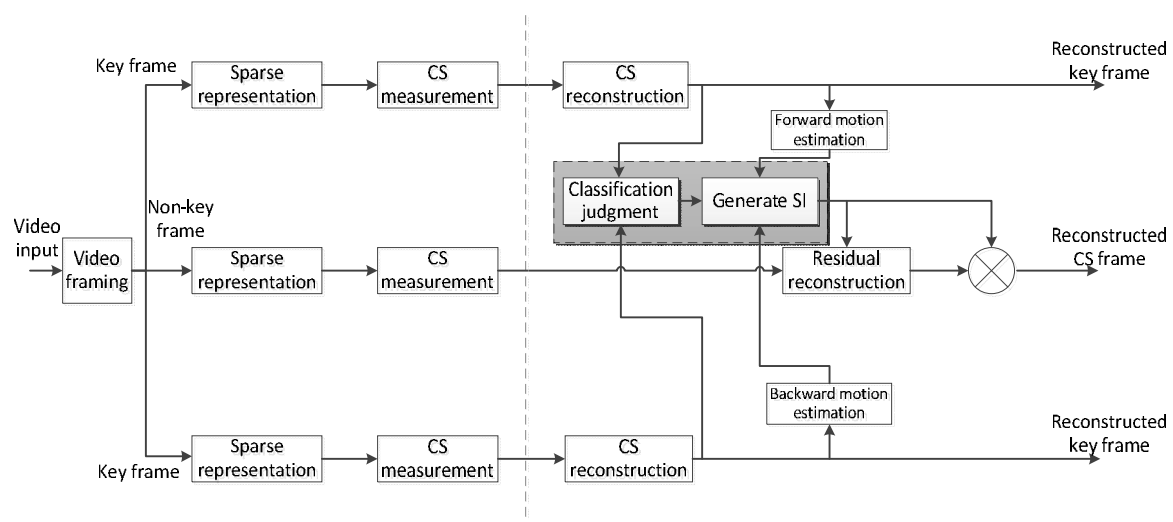


287 frames in the same scene are highly correlated. Therefore, we assume that even if the CS  
 288 measurement process is very different from a linear transformation (such as DCT), the CS  
 289 measurement of such adjacent frames is highly correlated. On the one hand, the DCT coefficients  
 290 follow a Laplacian distribution. On the other hand, CS measurements follow a more or less normal  
 291 (Gaussian) distribution.

## 292 4. DCVS Framework Based on Block Classification Weighted Side Information

### 293 4.1. Description of the Framework

294 In the DCVS framework, in order to improve the reconstruction quality of non-key frames at  
 295 the decoding end, side information is introduced when reconstructing non-key frames. Therefore,  
 296 the acquisition of side information plays a very important role in the DCVS framework. Especially  
 297 for the reconstruction of non-key frames, if the acquisition of side information is not accurate, the  
 298 performance of its reconstruction will be greatly affected. The process of obtaining the side  
 299 information may be generated by using a K-SVD training dictionary or performing motion  
 300 estimation on the decoded key frame, interpolating in the time domain to generate side information.  
 301 Most of the traditional methods of obtaining the side information are to estimate the forward and  
 302 backward motion of the reconstructed values of the two key frames before and after the current  
 303 frame, then according to a fixed weight (usually 1/2), adding to obtain the side information to assist  
 304 in the reconstruction of non-key frames [17]. However, the inter-frame correlations of different  
 305 regions in the video sequence are different, and the varying scenes and motion levels of the video  
 306 are different. If the inter-frame correlation between the forward and backward motion estimation is  
 307 weak, still using a fixed weight does not predict the current frame very well, so the accuracy of the  
 308 generated side information is not high, which in turn affects the reconstruction of non-key frames.  
 309 Therefore, the traditional method of synthesizing side information by using fixed weight 1/2 does  
 310 not make good use of inter-frame correlation. According to the correlation between different blocks  
 311 of video frames, a DCVS framework based on different blocks is proposed. The framework is  
 312 shown in Figure 3.



313

314 Figure 3. Block-based classification weighted side information DCVS framework

315 At the encoding end, the current frame and the two key frames are subjected to block sparse  
 316 representation and CS sampling measurement. The key frames are sampled by the traditional DWT  
 317 algorithm, the non-key frames are sampled by the hybrid sparse algorithm, and the measurement  
 318 matrix still selects the random Gauss matrix. At the decoding end, the key frame is first  
 319 reconstructed by SAMP algorithm, and then the classification judgment method introduced in the  
 320 previous section is used to classify non-key(CS) frames, performing forward and backward motion

321 estimation on reconstructed values of adjacent key frames, different weighting schemes are used to  
 322 generate side information according to different block categories. And the reconstruction of  
 323 non-key frames requires the use of side information in conjunction with residual decoding of  
 324 non-key frames.

#### 325 4.2. Non-key Frame Reconstruction

326 The reconstruction of non-key frames uses the motion estimation of the classification weights  
 327 to generate the side information, and then uses the side information and the measured values of the  
 328 current frame to reconstruct the residual of the current frame. The reconstruction result of the  
 329 non-key frame is the combination of the side information and the residual reconstruction.

330 Assuming that the current frame of a video is  $x$  and the predicted value of the current frame  
 331 is  $\tilde{x}$ , the residual between the actual result and the predicted value of the current frame is  
 332  $r = x - \tilde{x}$ , the predicted value  $\tilde{x}$  of the current frame is generated by predicting the current frame  
 333 from the pre- and post-decoding key frames. Since the decoding end has not yet obtained the  
 334 reconstructed data of the decoded key frames before and after the current frame, so we convert the  
 335 residual to the measurement domain:

$$336 \quad d = \Phi r = \Phi x - \Phi \tilde{x} = y - \Phi \tilde{x} \quad (8)$$

337 If the difference between the actual value and the predicted value of the current frame is  
 338 smaller, the smaller the residual of the two is, the sparser it is, the smaller the sampling under the  
 339 measurement matrix is, and the better the reconstruction effect of the residual is. The reconstruction  
 340 result of the current frame is:

$$341 \quad x_{rec} = \tilde{x} + r_{rec} \quad (9)$$

342 Where  $x_{rec}$  is the reconstruction result of the current frame,  $\tilde{x}$  is the predicted value of the current  
 343 frame, and  $r_{rec}$  is the result of the residual reconstruction.

344 In the acquisition of side information, it is necessary to utilize the reconstruction result of  
 345 adjacent decoding key frames [18]. First, forward and backward motion estimation is performed on  
 346 the reconstruction results of the  $i$ -th blocks  $x_{t-1}^i$  and  $x_{t+1}^i$  corresponding to the adjacent decoding  
 347 key frames to obtain  $\hat{x}_{t-1}^i$  and  $\hat{x}_{t+1}^i$  respectively, then find the side information of the  $i$ -th block  
 348 according to equation (10):

$$349 \quad SI^i = \alpha \hat{x}_{t-1}^i + (1 - \alpha) \hat{x}_{t+1}^i \quad (10)$$

350 Where  $\alpha$  is the weight coefficient.

351 In the previous section, the non-key frame blocks have been classified and judged. After the  
 352 blocks of different regions are classified, different weights  $\alpha$  are adopted to obtain the side  
 353 information, thereby further reconstructing the non-key frames. For the block with slight movement,  
 354 the ratio of the residual energy to the energy block of the previous key frame is small, and the  
 355 inter-frame correlation of the preceding and succeeding frames is large. According to equation (5),  
 356 the forward motion estimation can predict the current frame very well, and then take  $\alpha$  to a  
 357 larger value. Conversely, for the block with severe movement, the inter-frame correlation is small,  
 358 and the backward motion estimation can better predict the current frame, and then take  $\alpha$  to a  
 359 smaller value. For the block with big movement, the inter-frame correlation is moderate, and the  
 360 current frame can be predicted by weighted averaging of forward and backward motion estimation.  
 361 Based on the above analysis, the weights of the block with slight, big and severe movement are  
 362 respectively set to  $\alpha = \{0.7, 0.5, 0.3\}$ .

363 In summary, for the reconstruction process of non-key frames summarized as follows: Firstly,  
 364 the adjacent key frames are reconstructed, and the reconstruction results are classified and judged,  
 365 and the forward and backward motion estimation is performed on adjacent key frames. Combined  
 366 with the classification criteria and forward and backward motion estimation, the side information  
 367 of the non-key frame is obtained according to the equation (10). According to the motion estimation,  
 368 the side information and the current frame are obtained, and the residual reconstruction of the

369 non-key frame residual is performed. The reconstruction result of the non-key frame is the sum of  
 370 the side information obtained by the motion estimation and the non-key frame residual  
 371 reconstruction.

372 Algorithm 2: The description of the non-key frame reconstruction algorithm.

---

**Input:**  $y, \Phi_s, x_{t-1}, x_{t+1}$

---

1. Perform forward and backward motion estimation for the pre- and post-decoding key frames  $x_{t-1}$  and  $x_{t+1}$  respectively to obtain  $\hat{x}_{t-1}$  and  $\hat{x}_{t+1}$ ;
  2. According to the classification judgment result of  $x_{t-1}$  and  $x_{t+1}$ , the side information is obtained by the formula  $SI = \alpha x_{t-1} + (1 - \alpha)x_{t+1}$ , Where  $\alpha$  selects different weights depending on the classification;
  3. Calculate the residual of the measured value  $y$  and the SI in the measurement domain:  
 $r = y - cs\_Encoder(SI, \Phi_s)$ ;
  4. Reconstruct the residual  $r$  to  $\tilde{x}_r$ ;
  5. Reconstruction results for non-key frames:  $\tilde{x}_t = SI + \tilde{x}_r$ .
- 

**Output:**  $\tilde{x}_t$ : Reconstruction of the  $t$ -th block non-key frame.

---

### 373 5. Simulation Experiment Results and Analysis

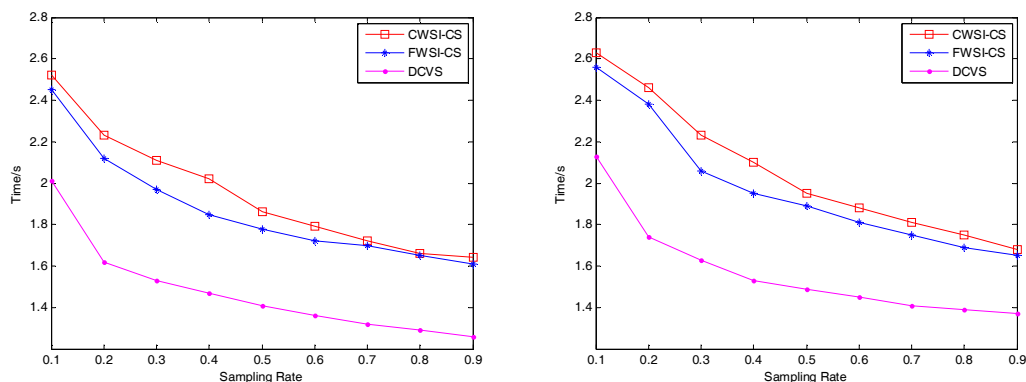
374 The experiment selects two sets of video sequences collected in the underground coal mine to  
 375 test the performance of the classified weighted side information method under hybrid sparse basis.  
 376 In the experiment, the classification weighted side information generation method and the fixed  
 377 weight side information generation method under the hybrid sparse basis are compared with the  
 378 original DCVS algorithm. The reconstruction algorithm in the classification weighted side  
 379 information selects the SAMP algorithm. The measurement and reconstruction of key frames and  
 380 non-key frames in the video sequence are block-based. We define the first frame as the key frame.  
 381 Since the key frame has a great influence on the generation and reconstruction of the side  
 382 information, so we choose a key frame with a sampling rate of 0.9. The non-key frame sampling  
 383 rate is selected by comparing the PSNR values reconstructed by the classified weighted side  
 384 information and the fixed weighted side information at different sampling rates. Select two  
 385 different video scenes, and the comparison results of PSNR values at different sampling rates are  
 386 shown in Figure 4. The classification thresholds  $T_1$  and  $T_2$  are:  $T_1=0.003$ ,  $T_2=0.015$ . In order to verify  
 387 the effect of the algorithm, by comparing the time complexity of the algorithm under different  
 388 sampling rates, the peak signal-to-noise ratio PSNR value is used to objectively evaluate the  
 389 reconstruction effect of the algorithm. To calculate the PSNR value, the value of the mean square  
 390 error MSE is first calculated. The MSE is defined as:

$$391 \quad MSE = \frac{1}{mn} \sum_{i=0}^{m-1} \sum_{j=0}^{n-1} \|I(i, j) - K(i, j)\|^2 \quad (11)$$

392 The peak signal to noise ratio PSNR is defined as:

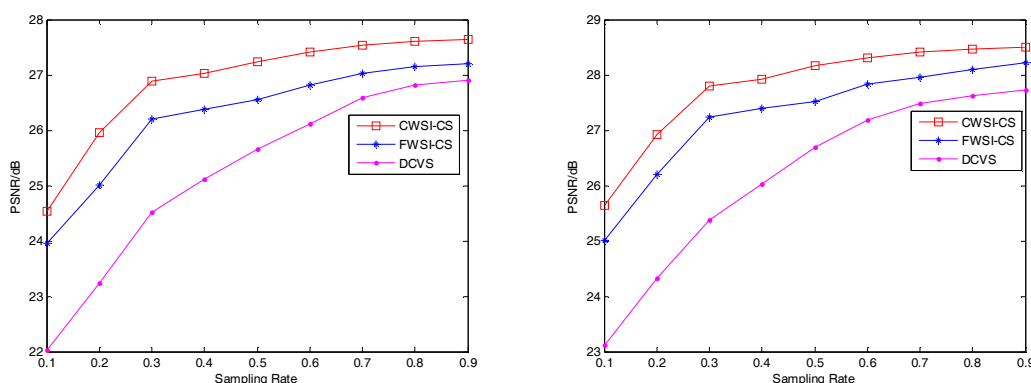
$$393 \quad PSNR = 10 \cdot \log_{10} \left( \frac{(2^n - 1)^2}{MSE} \right) \quad (12)$$

394 It is not difficult to understand that it can be obtained from equation (12) that the smaller the  
 395 value of the MSE is, the closer the reconstructed video frame effect is to the original video  
 396 frame image, and the larger the PSNR value is.



397  
398

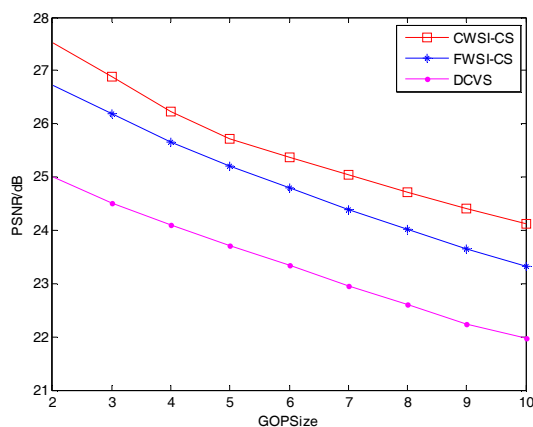
Figure 4. Relationship between sampling rate and time complexity of video sequences



399  
400

Figure 5. PSNR value of video reconstruction quality at different sampling rates

401 As can be seen from Figure 4, the complexity of the proposed algorithm is similar to the  
 402 fixed-weighted side information algorithm at different sampling rates. The original DCVS  
 403 algorithm has the shortest running time, that is, the complexity of the algorithm is lower than the  
 404 former two. However, as can be seen from Figure 5, as the sampling rate increases, the PSNR value  
 405 increases correspondingly. The reconstruction effect is relatively good under the condition of high  
 406 sampling rate. Obviously, the proposed algorithm has the highest reconstruction quality. Although  
 407 the complexity of its reconstruction is relatively high, its complexity is within an acceptable range  
 408 relative to the improvement of reconstruction quality. Figure 5 shows that when the sampling rate  
 409 is greater than 0.3, the increase of the PSNR value is not large as the sampling rate increases.  
 410 Considering that, in the case of meeting the general needs, the low sampling rate requires less  
 411 transmission data. Therefore, the sampling rate of non-key frames is selected to be 0.3, and the  
 412 lesser transmission data required at low sampling rates saves energy while ensuring reconstruction  
 413 quality.



414  
415

Figure 6. PSNR values of video reconstruction quality under different GOP sizes

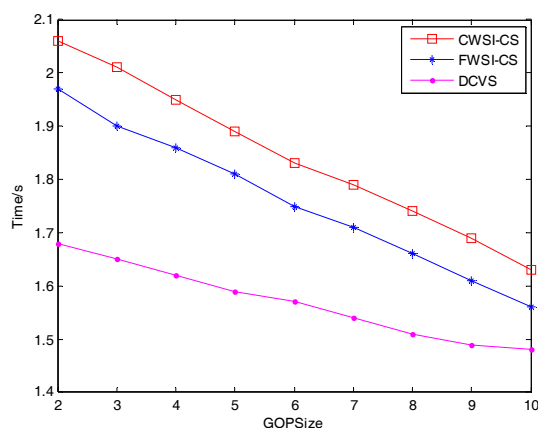


Figure 7. Time complexity under different GOP sizes

416  
417

418 As can be seen from Figure 6 and Figure 7, as the number of frames in each GOP increases, the  
419 inter-frame correlation becomes smaller, resulting in a decrease in the reconstructed PSNR value.  
420 However, the reconstruction effect of the proposed algorithm has always been above the fixed  
421 weight side information and the original DCVS algorithm, and its reconstruction complexity is  
422 close to the original DCVS algorithm, which saves energy consumption while ensuring the  
423 reconstruction quality.

424 Experiment to select the 30 GOPs of the video sequence, there are 3 frames in each GOP, and a  
425 total of 90 frames are respectively simulated. The two frames before and after the GOP are taken as  
426 key frames, and the intermediate frames are non-key frames. Through sparse representation, side  
427 information generation, frame reconstruction, the quality of the two side information generation  
428 method is compared between the classification weighted side information and the fixed weight side  
429 information, and is objectively measured by the PSNR value. The curve of the reconstructed quality  
430 PSNR value of 30 frames of non-key frames in the video sequence GOP is shown in Figure 8.

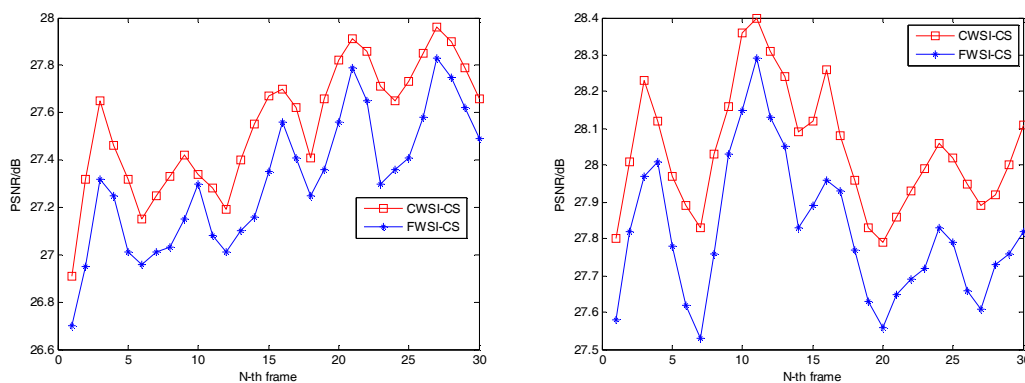


Figure 8. PSNR value of non-key frame reconstruction (sampling rate 0.3)

431  
432

433 It can be seen from the figure that under the same conditions, the PSNR value of the  
434 block-based classification adaptive sparse basis and the classification weighted side information is  
435 improved by 0.2-0.8dB compared with the fixed weight side information reconstruction. For image  
436 blocks with too high movement intensity, the correlation between adjacent frames is weak, and the  
437 reconstruction effect is relatively weak, but the reconstruction effect is still improved compared  
438 with the fixed weight side information. Therefore, the proposed algorithm can improve the quality  
439 of video reconstruction in different video scenarios.

## 440 6. Conclusions

441 In order to solve the problem that the correlation between the different video sequence image  
442 blocks is different in the process of applying the compressed sensing method to the video decoding  
443 side to acquire the side information for assisting non-key frame reconstruction, a block-based

444 classification weighting method is proposed for non-key frame reconstruction. In terms of the  
445 sparse representation of the encoding end, the video sequence is divided into different blocks  
446 according to the inter-frame correlation, and different sparse basis strategies are selected for  
447 different blocks; In the decoding end of the non-key frame, in the algorithm for generating the side  
448 information by motion estimation, the generation of the SI selects different weights according to the  
449 difference of the inter-frame correlation. The experimental results show that the scheme can  
450 adaptively select different sparse basis and side information generation schemes to assist non-key  
451 frame reconstruction according to different video scenes, making full use of inter-frame correlation  
452 and improving reconstruction quality.

453 **Author Contributions:** Conceptualization, K.W. and M.W.; Methodology, R.X.; Software, R.X.;  
454 Validation, R.X., K.W. and M.W.; Formal Analysis, R.X.; Resources, K.W.; Writing-Original Draft  
455 Preparation, R.X.; Writing-Review & Editing, H.X.; Visualization, R.X.; Supervision, H.X.; Project  
456 Administration, H.X.; Funding Acquisition, H.X.

457 **Funding:** This work is supported by the National Key Research and Development Project of China  
458 (No.2017YFC0804404).

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

## 460 References

- 461 1. Zhao, X.H.; Deng, Y.F.; Mu, D.C. Applied study on compressed sensing technology to mine Internet of  
462 things. *Coal Science and Technology*, **2016**, 44(7), 69-72 [[CrossRef](#)].
- 463 2. Zhao X.H.; Liu S.S.; Shen X.R.; et al. Micro-seismic data compression and reconstruction based on  
464 distributed compressed sensing. *Journal of China University of Mining & Technology*, **2018**, 1,  
465 172-182 [[CrossRef](#)].
- 466 3. Zhang F.; Yan X.X.; Li Y.J. A novel image reconstruction method of mine intelligent surveillance based on  
467 adaptive sparse representation. *Journal of China Coal Society*, **2017**, 42(5), 1346-1354 [[CrossRef](#)].
- 468 4. Zhang F.; Yan X.X. The block compressed sensing of mine monitoring images based on DFT basis. *Journal*  
469 *of Transduction Technology*, **2017**, 30(1), 94-100 [[CrossRef](#)].
- 470 5. Wang G.; Zhao Z.K.; Ning Y.J. Design of Compressed Sensing Algorithm for Coal Mine IoT Moving  
471 Measurement Data Based on a Multi-Hop Network and Total Variation. *SENSORS*, **2018**,  
472 18(6),1732 [[CrossRef](#)].
- 473 6. Lin B.L.; Zheng B.Y.; Qian C. The reconstruction methods of frames classification of distributed video  
474 compression coding. *Signal Processing*, **2015**, 2, 201-207 [[CrossRef](#)].
- 475 7. Zhao X.H.; Liu S.S.; Shen X.R.; et al. Research on processing algorithm of image in underground coal mine  
476 based on CS framework. *Coal Science and Technology*, **2018**, 46(2), 219-224 [[CrossRef](#)].
- 477 8. Zhang R. A preliminary study of image reconstruction and denoising based on compressed sensing.  
478 Southwest Jiaotong University, **2010** [[CrossRef](#)].
- 479 9. Wu M.H.; Zhu X.C. Dynamic measurement rate allocation for distributed compressive video sensing.  
480 *Journal of Nanjing University of Posts and Telecommunications*, **2013**, 33(1), 62-67 [[CrossRef](#)].
- 481 10. Dong H.; Zhuang B.; Su F.; et al. A novel distributed compressive video sensing based on hybrid sparse  
482 basis. // Visual Communications and Image Processing Conference. IEEE, **2014**, 320-323 [[CrossRef](#)].
- 483 11. Chen H.W.; Kang L.W.; Lu C.S. Dictionary learning-based distributed compressive video sensing.// *Picture*  
484 *Coding Symposium*. **2011**, 210-213 [[CrossRef](#)].
- 485 12. Aharon M.; Elad M.; Bruckstein A. rmK-SVD: An algorithm for designing overcomplete dictionaries for  
486 sparse representation. *IEEE Transactions on Signal Processing*, **2006**, 54(11), 4311-4322 [[CrossRef](#)].
- 487 13. Prades N.J.; Ma Y.; Huang T. Distributed video coding using compressive sampling.// *Picture Coding*  
488 *Symposium*. IEEE, **2009**, 1-4 [[CrossRef](#)].
- 489 14. Liu L.; Wang A.; Li Z.; et al. An Improved Distributed Compressive Video Sensing Based on Adaptive  
490 Sparse Basis. // *International Conference on Robot*. IEEE, **2011**, 137-140 [[CrossRef](#)].
- 491 15. Baig Y.; Lai E.M.K.; PUNCHIHEWA A. Distributed Video Coding Based on Compressed Sensing.//*IEEE*  
492 *International Conference on Multimedia and Expo Workshops*. IEEE, **2012**, 131(5), 325-330 [[CrossRef](#)].
- 493 16. Baraniuk R.; Davenport M.; Devore R.; et al. A Simple Proof of the Restricted Isometry Property for  
494 Random Matrices.// *CONSTR APPROX*, **2008**, 253-263 [[CrossRef](#)].

- 495 17. Dai Y.Y.; Cao X.Q.; Chen R.; et al. Reconstruction algorithm with classified weighted side Information for  
496 distributed video compressive sensing. *Computer Technology and Development*, **2017**, 27(5), 87-91[[CrossRef](#)].
- 497 18. Jian C.; Su K.X.; Wang W.X.; et al. Residual Distributed Compressive Video Sensing Based on Double Side  
498 Information. *Acta Automatica Sinica*, **2014**, 40(10), 2316-2323[[CrossRef](#)].