

***In silico* analysis on the functional and structural impact of Rad50 mutations involved in DNA strand break repair**

Juwairiah Remali¹, Wan Mohd Aizat², Chyan Leong Ng², Yi Chieh Lim³, Zeti-Azura Mohamed-Hussein^{2,4}, Shazrul Fazry^{Corresp. 1, 5}

¹ Department of Food Sciences, Faculty of Science and Technology, Universiti Kebangsaan Malaysia, Bangi, Selangor, Malaysia

² Institute of Systems Biology (INBIOSIS), Universiti Kebangsaan Malaysia, Bangi, Selangor, Malaysia

³ dKDanish Cancer Society, Research Centre Strand Boulevard, Copenhagen, Denmark

⁴ Department of Applied Physics, Faculty of Science and Technology, Universiti Kebangsaan Malaysia, Bangi, Selangor, Malaysia

⁵ Pusat Penyelidikan Tasik Chini, Fakulti Sains dan Teknologi, Universiti Kebangsaan Malaysia, Bangi, Selangor, Malaysia

Corresponding Author: Shazrul Fazry

Email address: shazrul@ukm.edu.my

Background. DNA double strand break repair is important to preserve the fidelity of our genetic makeup after DNA damage. Rad50 is one of the components in MRN complex important for DNA repair mechanism. Rad50 mutations can lead to microcephaly, mental retardation and growth retardation in human. However, Rad50 mutations in human and other organisms have never been gathered and heuristically compared for their deleterious effects. It is important to assess the conserved region in Rad50 and its homolog to identify vital mutations that can affect functions of the protein.

Method. In this study, Rad50 mutations were retrieved from SNPeffect 4.0 database and literature. Each of the mutations was analysed using various bioinformatic analyses such as PredictSNP, MutPred, SNPeffect 4.0, I-Mutant and MuPro to identify its impact on molecular mechanism, biological function and protein stability, respectively.

Results. We identified 103 mostly occurred mutations in the Rad50 protein domains and motifs, which only 42 mutations were classified as most deleterious. These mutations are mainly situated at the specific motifs such as Walker A, Q-loop, Walker B, D-loop and signature motif of the Rad50 protein. Some of these mutations were predicted to negatively affect several important functional sites that play important roles in DNA repair mechanism and cell cycle signalling pathway, highlighting Rad50 crucial role in this process. Interestingly, mutations located at non-conserved regions were predicted to have neutral/non-damaging effects, in contrast with previous experimental studies that showed deleterious effects. This suggests that software used in this study may have limitations in predicting mutations in non-conserved regions, implying further improvement in their algorithm is needed. In conclusion, this study reveals the priority of acid substitution associated with the genetic disorders. This finding highlights the vital roles of certain residues such as K42E, C681A/S, CC684R/S, S1202R, E1232Q and D1238N/A located in Rad50 conserved regions, which can be considered for a more targeted future studies.

***In silico* analysis on the functional and structural impact of Rad50 mutations involved in DNA strand break repair**

Juwairiah Remali¹, Wan Mohd Aizat², Chyan Leong Ng², Yi Chieh Lim³, Zeti-Azura Mohamed-Hussein^{2,4}, and Shazrul Fazry^{1,5,*}

¹ Department of Food Sciences, Faculty of Science and Technology, Universiti Kebangsaan Malaysia, Bangi, Selangor, Malaysia

² Institute of Systems Biology (INBIOSIS), Universiti Kebangsaan Malaysia, Bangi, Selangor, Malaysia

³ dKDanish Cancer Society, Research Centre Strand Boulevard, Copenhagen, Denmark

⁴ Department of Applied Physics, Faculty of Science and Technology, Universiti Kebangsaan Malaysia, Bangi, Selangor, Malaysia

⁵ Pusat Penyelidikan Tasik Chini, Fakulti Sains dan Teknologi, Universiti Kebangsaan Malaysia, Bangi, Selangor, Malaysia

Corresponding Authors:

Shazrul Fazry

Department of Food Sciences, Faculty of Science and Technology, Universiti Kebangsaan Malaysia, Bangi, Selangor, Malaysia Email address: shazrul@ukm.edu.my

Abstract

Background. DNA double strand break repair is important to preserve the fidelity of our genetic makeup after DNA damage. Rad50 is one of the components in MRN complex important for DNA repair mechanism. Rad50 mutations can lead to microcephaly, mental retardation and growth retardation in human. However, Rad50 mutations in human and other organisms have never been gathered and heuristically compared for their deleterious effects. It is important to assess the conserved region in Rad50 and its homolog to identify vital mutations that can affect functions of the protein.

Method. In this study, Rad50 mutations were retrieved from SNPeffect 4.0 database and literature.

Each of the mutations was analysed using various bioinformatic analyses such as PredictSNP, MutPred, SnpEff 4.0, I-Mutant and MuPro to identify its impact on molecular mechanism, biological function and protein stability, respectively.

Results. We identified 103 mostly occurred mutations in the Rad50 protein domains and motifs, which only 42 mutations were classified as most deleterious. These mutations are mainly situated at the specific motifs such as Walker A, Q-loop, Walker B, D-loop and signature motif of the Rad50 protein. Some of these mutations were predicted to negatively affect several important functional sites that play important roles in DNA repair mechanism and cell cycle signalling pathway, highlighting Rad50 crucial role in this process. Interestingly, mutations located at non-conserved regions were predicted to have neutral/non-damaging effects, in contrast with previous experimental studies that showed deleterious effects. This suggests that software used in this study may have limitations in predicting mutations in non-conserved regions, implying further improvement in their algorithm is needed. In conclusion, this study reveals the priority of acid substitution associated with the genetic disorders. This finding highlights the vital roles of certain residues such as K42E, C681A/S, CC684R/S, S1202R, E1232Q and D1238N/A located in Rad50 conserved regions, which can be considered for a more targeted future studies.

Keywords: Rad50, Rad50 related diseases, DNA damage, Rad50 mutation, Rad50 *in silico* model

Introduction

DNA repair process exists in all organisms including both prokaryotes and eukaryotes, and most of the related proteins in this process are known to be highly conserved throughout biological evolution. One of such protein complexes involved in eukaryotic DNA repair process is MRN complex and it comprises of three proteins i.e. meiotic recombination 11 (Mre11), DNA repair protein Rad50, and nibrin (called Nbn or Nbs1). These proteins play an important role in maintaining the genomic integrity by orchestrating DNA damage checkpoint, telomere maintenance, homologous recombination (HR) as well as non-homologous end joining repair (NHEJ) mechanism (van den Bosch, Bree & Lowndes, 2003). MRN complex is one of the first factors to be localised to DNA lesions where it has a structural role by tethering and stabilising broken chromosomes (de Jager et al., 2001; van den Bosch, Bree & Lowndes, 2003).

Null mutations in MRN complex have been shown to be lethal in higher eukaryotes such as in embryonic stem cells (Luo et al., 1999). In addition, mutations in the *Nbs1* gene, can cause Nijmegen breakage syndrome (NBS) whereas Mre11 mutations resulted in Ataxia telangiectasia-like disease syndrome (ATLD) (Carney *et al.*, 1998). So far, studies of Nbs1 and Mre11 deficiencies in human have been extensively investigated through cells and clinical data obtained

from NBS and ATLD patients (Barbi et al., 1991; Waltes et al., 2009). Unfortunately investigation of the effect of Rad50 mutations on human is very limited due to only one patient with fully characterized Rad50 deficiency (known as NBS like disorder (NBSLD)) has been reported (Waltes et al., 2009). This NBSLD patient with microcephaly, bird-like features, radiosensitivity and delayed development, revealed to inherit heterozygous mutation from her parents (Barbi et al., 1991). The first mutation (c.3277C/T; p.R1093X) on exon 21 was maternally inherited causing a premature termination codon, thus producing a truncated Rad50 protein. Whereas the second mutation on the exon 25 (c.3939A/T) was paternally inherited and it has changed the stop codon of normal Rad50 to a tyrosine codon, thereby producing a larger Rad50 protein (Waltes et al., 2009). Both mutation interestingly give rise to hypomorphic characterization of the Rad50 expressions in this patient (Gatei et al., 2011). The cause of this characteristic are still on debate to this day. Given that perturbation of Rad50 structure and function could contribute to genomic instability (Assenmacher & Hopfner, 2004), it is therefore important to decipher its conserved domains and genetic polymorphism.

Single nucleotide polymorphism (SNP) is one of the most common types of genetic variation in human (Lee et al., 2005). Even though most of the polymorphic changes do not affect normal cellular function, some variants do influence gene expression or translated protein function (Risch & Merikangas, 1996; Collins, Guyer & Charkravarti, 1997). For instance, cystic fibrosis (Bartoszewski et al., 2010), sickle-cell anemia (Shaikho et al., 2017), and β -thalassemia (Traeger et al., 1980) are examples of diseases resulted from SNPs. Nearly half of the disease-related mutations are derived from nonsynonymous SNPs (nsSNPs), a single base change that alters the amino acid sequence of the encoded protein (Cargill et al., 1999; Halushka et al., 1999). Although it is remarkably important to reveal the connection between SNPs and related diseases, the accelerating number of known SNPs have made it very difficult to discriminate between pathogenic and neutral variants through experimental validations (Tranchevent et al., 2011). Therefore, bioinformatic prediction tools have become extremely critical for the initial analysis of their molecular functions as well as prioritization of further experimental characterization including deciphering the effects of Rad50 SNPs (Bendl et al., 2014). Furthermore, prioritization of disease candidates genes from experiment and databases evidence is essential for further pathological investigation (Piro & Di Cunto, 2012). Several investigations on Rad50 mutations have been reported in human (Waltes et al., 2009; Gatei et al., 2011), mice (Bender et al., 2002; Roset et al., 2014), yeast (Alani, Padmore & Kleckner, 1990; Chen et al., 2005), and archaea (Koroleva et al., 2007) yet there are still no reports that compare these experimental results with *in silico* prediction, which will be important for the protein functional annotation. Moreover, a number of different SNPs for Rad50 have been deposited in SNP databases but their impact on the cellular regulation have not been thoroughly investigated thus far.

Hence, the aim of this study was to identify the functional and structural effects of amino acid mutations in Rad50 gathered from exhaustive literature review and SNP database (SNPeffect

4.0) search. Rad50 sequences in different organisms including human and selected animals (chimpanzee, rats, mice, zebra fish, rabbit and fruit fly) and yeasts were compared and aligned to identify their conserved residues. Mutations that contributed to the most damaging effects were then analysed *in silico* using PredictSNP for the amino acid impact after the substitution, MutPred for predicting molecular mechanism, SNPeffect for identification of protein or amyloid aggregation as well as I-Mutant and MuPro for protein stability after the mutation. Such approach was also successfully reported by several researchers studying the impact of various SNPs. For example, Marín-Martín *et al.*, (2014) studied the impact of SNPs in the ABCA1 transporter gene by cross validating their prediction with experimentally reported data. Another study by Fawzy *et al.*, (2015) also validated their *in silico* approach finding by means of comparison with available literature to study gene polymorphisms in obese children and adolescents. In this study, Rad50 mutations gathered from various studies are compared with their *in silico* predictions. This is highly valuable in understanding Rad50 functional roles especially during DNA strand break, allowing prioritization of mutations or sites to be studied in future *in vivo* studies, whilst bearing in mind its possible impact on human. Ultimately, this may help on the development of precision medicine for Rad50 mutations in human.

Materials & Methods

Multiple sequence alignment (MSA) analysis and conserved domain analyses

Human Rad50 protein sequence was obtained from National Center for Biotechnology Information (NCBI). The sequence similarity search tool, BLASTP from the NCBI server (<http://blast.ncbi.nlm.nih.gov/Blast.cgi>) was used to find homologs for Rad50. To investigate the similarity between Rad50 protein in human and other organisms such as *Danio rerio*, *Mus musculus*, *Rattus norvegicus*, *Pan troglodytes*, *Oryctolagus cuniculus*, *Drosophila melanogaster*, *Saccharomyces cerevisiae* and *Schizosaccharomyces pombe*, a multiple sequence analysis (MSA) was conducted using Clustal Omega (<https://www.ebi.ac.uk/Tools/msa/clustalo/>) with default settings to determine consensus and conserved regions between the multiple sequences of different organisms (Sievers & Higgins, 2018). Meanwhile, InterPro (<http://www.ebi.ac.uk/interpro/>) was used to identify the domains and motifs using human sequence (Finn et al., 2017). InterPro results are classified into several types (families, domains, motif or sites) depending on the biological entity they represent (Finn et al., 2017). Using this tool, Rad50 protein sequence was classified into families and the presence of domains and important sites were predicted. ClustalX software (Thompson, Gibson & Higgins, 2002) was used to view and analyse the conserved regions within the domains and motifs in the selected proteins.

Data mining of Rad50 mutation from literature and SNPs database

Rad50 mutations were identified from previous published manuscripts using PubMed database and their functional impacts were extracted for comparison. Besides that, naturally occurring single nucleotide polymorphisms (SNPs) in Rad50 were retrieved from SNPeffect 4.0 database (<http://snpeffect.switchlab.org/about>) (De Baets et al., 2012) (date of access: April 07th 2018). SNPeffect 4.0 database currently contains more than 60,000 human SNPs gathered from human avariance list available at UniProt website (<https://www.uniprot.org/>). It specifically focuses on the molecular characterization, annotation of diseases as well as polymorphism variants in human proteins (De Baets et al., 2012). All these available Rad50 protein mutations (obtained from both literature and databases) have been aligned using pairwise alignment through Clustal Omega between human sequence and other organisms' sequence, individually. From this analysis, we identified similar mutation sites in human. All the identified equivalent mutations in human were manually refined, for example removing the same residues and mutations that has been studied by several different researchers to identify the non-redundant mutations (S1 Table). Identified mutations (after converting to equivalent residues in human) were then mapped into Fig. S1.

Secondary structure prediction and analysis of 3D modeling

The Rad50 templates identified from the BLAST analysis also were used to develop secondary structure and 3D model. The PSIPRED program (<http://bioinf.cs.ucl.ac.uk/psipred/>) has been utilised for secondary protein structure prediction (Buchan et al., 2013). Secondary structure prediction has revealed a clear distribution of alpha helix, beta sheet and coil in *H. sapiens* (Helix: 74.69%, coil; 18.29 and beta sheet; 7.01% (S2 Figure). Databases such as UniProt (<https://www.uniprot.org/>) and Protein Data Bank (PDB) (<https://www.rcsb.org/>) were used to identify structural information regarding Rad50 protein in human. Rad50 protein sequence also has been BLAST searched against Protein Data Bank (PDB) sequence in Network Protein Sequence @analysis (NPS@) (<https://npsa-prabi.ibcp.fr/>) to identify the most identical structure. The incomplete structure has been further predicted using fold recognition method using Protein Homology/analogY Recognition Engine Version 2.0 (Phyre2) (<http://www.sbg.bio.ic.ac.uk/phyre2>) (Kelley et al., 2015). Phyre2 is an online tool to predict and analyze protein structure, function and mutations which uses advanced remote homology detection methods to build 3D models, predict ligand binding sites and analyze the effect of amino acid variants (e.g., nonsynonymous SNPs (nsSNPs)) for a protein sequence (Kelley et al., 2015). Rad50 sequence was submitted to the webserver to interpret the secondary and tertiary structures of the model, domain composition and quality. 3D model of Rad50 was run under 'intensive' mode that generates a complete full-length model of a protein sequence by using multiple template modelling and simplified *ab initio* folding simulation (Kelley et al., 2015).

UCSF Chimera software was used to view and to analyse the 3D structure (Pettersen et al., 2004).

Prediction of deleterious effects of Rad50 mutations using *in silico* tools

The Rad50 mutations were *in silico* predicted using PredictSNP to determine their possible molecular impacts in human (<https://loschmidt.chemi.muni.cz/predictsnp1/>) (Bendl et al., 2014). Its benchmark dataset contains over 43,000 mutations obtained from the Protein Mutant Database and the UniProt database (Bendl et al., 2014). This tool incorporated six established prediction tools; such as Multivariate Analysis of Protein Polymorphism (MAPP) (Stone & Sidow, 2005), Predictor of human Deleterious Single Nucleotide Polymorphisms (PhD-SNP) (Capriotti & Fariselli, 2017), PolyPhen-1 (Ramensky, 2002), PolyPhen-2 (Adzhubei, Jordan & Sunyaev, 2013), Sorting Intolerant from Tolerant (SIFT) (Sim et al., 2012) and Single-Nucleotide Amplified Polymorphisms (SNAP) (Bromberg & Rost, 2007) to provide a more accurate and robust comparison. We classified the mutations as deleterious if five to seven analyses performed were identified as damaging in PredictSNP. For instance, an *in silico* prediction was considered accurate when a given mutation predicted to be deleterious (as performed in this study) was also found experimentally deleterious (either *in vitro* or *in vivo* with phenotypes such as embryonic lethality, growth defect and/or cancer predisposition) based on previous cited studies. Conversely, the prediction is inaccurate if such deleterious mutations was predicted as neutral or tolerant.

Molecular mechanism of amino acid substitutions

To determine the molecular mechanism based on pathogenicity of amino acid substitutions in Rad50, MutPred2 (Pejaver et al., 2017) (<http://mutpred2.mutdb.org/index.html>) analysis was carried out. This program predicts the pathogenicity and molecular impacts of amino acid substitutions potentially affecting the phenotype. It is trained on a set of 53,180 pathogenic and 206,946 unlabelled (putatively neutral) variants obtained from the Human Gene Mutation Database (HGMD) (Stenson et al., 2017), SwissVar (Mottaz et al., 2010), dbSNP (Sherry, 2001) and inter-species pairwise alignment (Pejaver et al., 2017). The output of MutPred contains a general probability that the amino acid substitution is deleterious/disease-associated, and a list of rank of specific molecular alterations potentially affecting the phenotype with its *p*-value (<0.05).

Prediction of molecular and structural effects of protein coding variants in Rad50 mutation

Prediction of molecular and structural effects of protein coding variants in Rad50 mutations was performed using SnpEff4.0 (De Baets et al., 2012) (<http://snpeff.switchlab.org/about>). The analysis includes predictions of the aggregation prone regions in a protein sequence (TANGO),

amyloid-forming regions (WALTZ) and chaperone binding site (LIMBO). The range of prediction score differences outside -50 to 50 for mutants are considered significant (De Baets et al., 2012). SNPEffect also uses FoldX (Schymkowitz et al., 2005) to analyse the effect of mutations on the structural stability. However, as structure quality is important for the accuracy of delta G predictions for stability, model structures with less than 90% sequence identity to the modelling template structure will not be modelled (De Baets et al., 2012).

Analysis of protein stability

The stability of Rad50 upon single amino acid residue mutations were predicted using MUpro (<http://mupro.proteomics.ics.uci.edu/>) (Cheng, Randall & Baldi, 2006) and I-Mutant 3.0 (<http://gpcr2.biocomp.unibo.it/cgi/predictors/I-Mutant3.0/I-Mutant3.0.cgi>) (Capriotti, Fariselli & Casadio, 2005) using default setting, for instance temperature was set at 25°C and pH 7. MUpro and I-Mutant 3.0 are valuable tools for protein stability prediction and analysis, even when the protein structure is not yet known with atomic resolution. Both use support vector machines (SVM)-based tool to predict protein stability changes for single amino acid mutations either from both sequence or structural information which correctly predicts with over 80% accuracy using cross validation methods (datasets and experimental) (Capriotti, Fariselli & Casadio, 2005; Cheng, Randall & Baldi, 2006). Rad50 protein sequence was searched against the web server and energy changes ($\Delta\Delta G$) were recorded. Negative value for $\Delta\Delta G$ represents a decrease in protein stability whereas positive value for $\Delta\Delta G$ represents an increase in stability.

Results

Rad50 data acquisition and MSA analysis

Human Rad50 sequence from NCBI database contains 1312aa with the accession number of AAB07119.1. Sequence homology search of the human Rad50 protein was performed against NCBI nonredundant protein databases (E -value $\leq 1E-05$) and the result was downloaded for further analysis. Out of 500 sequences, six sequences were chosen for MSA analysis from diverse organisms such as *D. rerio*, *M. musculus*, *R. norvegicus*, *P. troglodytes*, *O. cuniculus*, and *D. melanogaster*. Two sequences, *Saccharomyces cerevisiae* and *Schizosaccharomyces pombe* were also included due to widely being used as models in previous Rad50 studies (S1 Table).

Analysis of protein domains

Domain identification analysis showed that Rad50 contains three P-loop containing nucleoside triphosphate hydrolase (P-loop NTPase) domains which belong to ATP Binding Cassette (ABC) protein superfamily (De La Rosa Metzere Bierlein & and Scott W. Nelson, 2011). It is located

near the N- and C-terminal, at the residue number of 25-103, 130-227 and 1196-1279 (Figure 1a). Residue annotation showed that Rad50 has six specific motifs including Walker A and Q-loop that are located at the N-terminal whereas Rad50 signature motif, Walker B, D-loop and H-loop/switch region are located at C terminal (Figure 1a) (De La Rosa Metzere Bierlein & Scott W. Nelson, 2011). It also has a domain called zinc hook (635-734aa) located at C-terminal region (Figure 1a) (Hopfner et al., 2002). Multiple sequence alignment (MSA) analysis between human Rad50 and its homologous genes (*D. rerio*, *M. musculus*, *R. norvegicus*, *P. troglodytes*, *O. cuniculus*, *D. melanogaster*, *S. cerevisiae* and *S. pombe*) also revealed that these specific motifs are highly conserved (Figure 1b).

Mutation datasets from the literature and database searches

In order to identify the Rad50 mutations, literature pertaining to the topic was exhaustively searched and 18 articles over the period of 1990 to 2017 were identified. All these mutations from different organisms were listed in S1 Table. There are 103 mutations identified which mostly occurred in the protein domains and motifs with various biological effects (S2 Table). In order to obtain equivalent mutations in human, pairwise alignment was performed individually between each organism (*D. rerio*, *M. musculus*, *R. norvegicus*, *P. troglodytes*, *O. cuniculus*, *D. melanogaster*, *S. cerevisiae* and *S. pombe*) and the Rad50 human sequence as a reference (S1 Table). Then, MSA analysis was carried out between these sequences from different organisms (including human) to identify consensus regions (Figure 1 and S1). Further refinement such as integrating similar mutations that occurred at the same positions (for examples; S1202R, K42R, S679R, P682E, V683R, R1214E, K6E, and R81I (S2 Table) from different organisms of which a total of 80 different mutations or non-redundant mutation were identified. All these mutations have been mapped based on equivalent residues in human (S1 Figure). From SNPeffect 4.0 database, another 13 SNP mutations were also identified (S3 Table). However, from the total of 103 mutations obtained from literature, only 42 residues of the Rad50 protein mutations were known to contribute to the most damaging effects *in vitro* and *in vivo* such as embryonic lethality (Bender et al., 2002; Roset et al., 2014) and growth defect (Table 1 and S2 Table) (Alani, Padmore & Kleckner, 1990; Bhaskara et al., 2007; Waltes et al., 2009; He et al., 2012; Barfoot et al., 2015; Hohl et al., 2015). Most of these deleterious mutations reside at the specific motifs such as Walker A, Q-loop, zinc hook, Rad50 signature motif, Walker B and D-loop (Figure 1b) that become our primary research focus (Figure 1b).

3D structure modelling of Rad50

Currently there is no complete structure of human Rad50 available. Nonetheless, a crystal structure of Rad50 hook and coil-coil domain (HCC) that contains 182 residues (residue 585-766) has been determined (PDB ID: 5GOX) (Park et al., 2017), which represents 13% of the

Rad50 structure in human. We attempted to predict a more complete human Rad50 3D structure model using homology modeling. Homology modeling program Phyre2 successfully predict the N-terminal of 276 residues (2-278) and C-terminal of 155 residues (1153-1306) of human Rad50 with 100% confidence level using the template model (PDB ID: 5DAC) from *Chaetomium thermophilum* (Seifert et al., 2016) that share 66% sequence identity (S2 Figure). As a result, half of human Rad50 protein structure was obtained. The regions with no 3D structure information available are residues 279 to 584 and residues 767 to 1152 (S2 Figure), which were mainly predicted to consist of alpha helices secondary structure (S2 Figure). Pairwise alignment of Rad50 sequence between *C. thermophilum* (1315aa) and human (1312aa) showed about 30% sequence identity (S2 Figure). The result showed that the partial structure of *C. thermophilum* that has been determined (black line) are highly conserved with the human sequence (S2 Figure) suggesting that the human structure should also share high structure similarity to *C. thermophilum* at these regions. In agreement to this, results from Phyre2 prediction showed that the N-terminal and C-terminal of Rad50 form a globular and coil-coil domain, similar to the structure of *C. thermophilum* (Figure 2a). With the generated model, six motifs of Rad50 namely Q-loop, Walker B, signature motif, D-loop, Walker A and H-loop were identified and marked in the 3D structure (Figure 2a). All identified deleterious residues found in the domain were also marked as shown in Figure 2b. To correlate the deleterious residues in the Rad50 HCC domain with zinc hook motif that was not found in the model, the structure of 181 residues (residue 585-766) that has been determined (PDB ID: 5GOX) independently was employed (Figure 2c and 2d) for functional analysis.

Analyses of Rad50 mutation deleterious effects

All of these 42 mutations (based on mutations from other organisms mapped to human) (Figure 1b and 2b) were then analysed using bioinformatics analyses such as impact of amino acid substitutions (PredictSNP), molecular mechanism (MutPred), structural phenotyping (protein and amyloid aggregation) (SNPeffect 4.0) and protein stability (MuPro and I-Mutant 3.0) (Table 2). All raw data from each analysis has been supplied as supplementary data (S3 Table for PredictSNP analysis, S4 Table for MutPred analysis, S5 Table for SNPeffect analysis and finally S6 Table for I-Mutant and MuPro analysis). The results showed that most of the deleterious effects fall into specific motifs such as Walker A, Q-loop, Rad50 signature motif, Walker B and D-loop (Table 2 and Figure 2a). Previous analysis also revealed that mutations at these motifs contributed to a number of biological defects such as growth defect (Alani, Padmore & Kleckner, 1990; He et al., 2012; Hohl et al., 2015) embryonic lethality (Bender et al., 2002; Roset et al., 2014), cancer predisposition (Bender et al., 2002; Roset et al., 2014), hematopoietic and spermatogenic depletion (Bender et al., 2002; Roset et al., 2014) (Table 1). Several mutations at the zinc hook region (C681A, C681S, P682E, C684R and C684S) and ATPase/coiled-coil domain (K6E and K132E) also showed to be deleterious (Table 2, Figure 2c and 2d).

Furthermore, mutations located at Walker A (Figure 2a) were predicted to affect catalytic and allosteric site, loss or gain of methylation, alteration of DNA binding, metal binding, ordered interface and the loss of relative solvent accessibility, which all are depending on the types of amino acid substitutions (Table 2). These mutations were predicted to affect ATP binding site motif, N-myristoylation, casein kinase II (CK2), protein kinase A (PKA) phosphorylation site and Forkhead-associated (FHA) functional sites. Mutations at the Walker A region also might led to the decrement of protein stability as predicted by I-Mutant and MuPro. Mutation at the Q-loop region (Q159H) (Figure 2a and 2b) also predicted to have significant deleterious effect and decreased protein stability, but no effects have been identified on its molecular mechanism and structural phenotyping as predicted by MutPred and SNPeffect 4.0, respectively (Table 2).

Mild deleterious effect was predicted at the mutated zinc hook domain (Table 2 and Figure 2e). Subsequent analysis using MutPred also revealed that any mutation at zinc hook might affect several important functional sites that involved in DNA damage repair signalling response and cell cycle checkpoints such as phosphatidylinositol 3-kinase-related kinases (PIKK) phosphorylation site, protein kinase C (PKC) phosphorylation site and BRCA1 C-terminus (BRCT) phosphopeptide ligands binding sites (Table 2). Moreover, deleterious mutation was predicted at the conserved cysteine residue located at the zinc hook motif (CXXC). For example, amino acid substitutions of alanine (A) and serine (S) at the cysteine residue position 681; (C681(A/S)) (Figure 2d) may affect N-glycosylation, proline-directed phosphorylation and mitogen-activated protein kinases (MAPK) phosphorylation site, which possibly due to the affected zinc binding domain (Table 2). Another deleterious mutation, C684(R/S) was also predicted to not affect its molecular mechanism but might disrupt diarginine retention/retrieving signal, PKC and PIKK phosphorylation site (Table 2). Whilst P682E (Figure 2d) mutation may lead to gain of helix, altered coiled coil domain, loss of N-linked glycosylation and CK2 phosphorylation site (Table 2).

Rad50 signature motif (Figure 2a) is a critical site which could lead to deleterious effects if mutated as suggested by PredictSNP analysis (Table 2). All mutations in this motif (S1202A/R/M, Q1205E and K1206M/A/E) or located near this motif (G1198E, L1211W and R1214A/E) (Figure 1b and -2b) were predicted to affect the protein allosteric and catalytic sites (Table 2), except for R1198E. Mutations at residue S1202A/R/M (Figure 1b and 2b) might affect PKA phosphorylation sites and glycosaminoglycan attachment site (Table 2). Furthermore, R1214A (Figure 1b and 2b) mutation might affect ATP-binding cassette, ABC transporter-type, signature and profile functional sites (Table 2). We have also predicted ~~identified~~ several mutations in Rad50 signature motif such as Q1205E, L1211W and R1214A that contributed to the total defect in the structural phenotyping such as the increment in protein and amyloid aggregation and the decrement of protein stability (Table 2).

We have also predicted K6E, K132E and K105E mutations occurred at the coil-coiled domain or ATPase domain to be deleterious (Figure 1b, 2b and Table 2). Specifically, the mutations at K6E and K132E might lead to loss of strand or loss of helix, respectively. Additionally mutation at K132E also predicted to affect casein kinase 1 (CK1) and PKC phosphorylation sites (Table 2). Even though K22M and R83I (Figure 2b) were predicted to be neutral in PredictSNP analysis, both of these mutations have also been predicted to increase protein aggregation tendency (Table 2). The mutation at R83I might contributed to the alteration of coiled coil structure domain, DNA binding and ordered interface, that might affect the functional site such as protein-protein interactions (PPI)-docking motif (Table 2). Another neutral mutation predicted were T191E, C221E and S106E (Figure 1b and 2b), where T191E mutation might be responsible in altering the coiled coil domain and may affect tumor necrosis factor receptor-associated factor (TRAF), serine/threonine-protein kinase (NEK2) and PKC phosphorylation site (Table 2). On the other hand, C221E and S106E (Figure 1b and 2b) were predicted to not affect any molecular mechanism or protein aggregation (Table 2).

Discussion

Rad50 is a member of the structural maintenance of chromosomes (SMC) family of proteins that participates in chromosome structural changes (Kinoshita et al., 2009). The globular ABC ATPase head domain is formed by the N- and C-termini (Figure 2a) (Hohl et al., 2011). The coiled-coil apex of Rad50 contains a conserved cysteine amino acid motif across the organisms, which is called the zinc hook (Kinoshita et al., 2009). When DNA double strand break occurs, Rad50 complex binds to the DNA early in the repair process to recognize such breaks and grips them in close juxtaposition (Paull & Gellert, 1998; de Jager et al., 2001). This protein also activates ATM kinase that is crucial for DNA damage signalling (Uziel et al., 2003).

Rad50 globular head domain contains conserved domains and motifs (Figure 1a and 2a) such as P-loop NTPase domains and six motifs which are Walker A and B motifs, Rad50 signature motif, D-loop, H-loop, and Q-loop motif (Figure 1a and 2a). P-loop NTPase domains are belong to ABC protein superfamily. The ABC protein superfamily has been identified in diverse organisms and is also known to be one of the most conserved protein superfamilies (Jones, O'Mara & George, 2009). ABC proteins consist of six conserved motifs (Figure 1a and 2a) which make up the nucleotide binding domain in Rad50 (Symington, 2002). The nucleotide binding domain of ABC protein is known to play an important role in binding and hydrolyzing ATP at its dimeric interface (Davidson et al., 2008). Rad50 also has a special conserved Cys-X-X-Cys zinc hook motif at the center of coiled-coil domain (Figure 2c). This motif is in a hook-shaped structure which dimerizes a second hook via cysteine-mediated zinc ion coordination (Figure 2c) (Hopfner et al., 2002). This zinc dependent dimerization event allows the formation of MRN complex which has suitable lengths and conformational arrangements to link sister chromatids in HR and DNA ends in NHEJ (Hopfner et al., 2002).

Consistency between bioinformatics prediction and experimental evidence

PredictSNP was used in this study to provide a more accurate prediction of disease-related mutations as it combines six best performing prediction tools for a consensus classifier (Bendl et al., 2014). Evidently, this *in silico* analysis was consistent with the results from the previous experimental studies where mutations at the Walker A, D-loop, signature motif, Q-loop and Walker B have shown damaging effects (Table 1, 2 and Figure 2a).

G41D and K40E (Figure 1b and 2b) mutations at the Walker A motif (Figure 2a) and C681A and C684R (Figure 1b and 2d) mutations at the cysteine residue (CXXC) in the zinc hook motif (Figure 2c) conferred an identical phenotype with the Rad50 null mutation characterized by total defect in the formation of viable spore in *S. cerevisiae* experiment (Table 1) (Alani, Padmore & Kleckner, 1990; He et al., 2012). This analysis also identified that mutations at Q-loop (Q159H) and D-loop (D1238N and D1238A) (Figure 1b, 2a and 2b) were also predicted deleterious (Table 2) and were experimentally shown to interrupt all ATP-dependent activities of the complex in different organisms such as *P. furiosus* and T bacteriophage respectively (Table 1) (Moncalian et al., 2004a; De La Rosa Metzere Bierlein & and Scott W. Nelson, 2011). Furthermore, a E1232Q (Figure 1b and 2b) mutation at the Walker B motif (Figure 2a) was also predicted to be deleterious (Table 2). Similarly the mutation of Walker B at residue E798Q in *Thermotoga maritima* showed low ability to respond to DNA damage (Table 1) (Rojowska et al., 2014a). This suggests that this motif is important for a molecular repair process, specifically during DNA binding process, which if mutated will affect the viability of an organism. Our analysis using PredictSNP has identified three mutations, which were N28A (Figure 1b and 2b) (De La Rosa Metzere Bierlein & and Scott W. Nelson, 2011), D1238H (Figure 1b and 2b) (De La Rosa Metzere Bierlein & and Scott W. Nelson, 2011) and S1202R (Figure 1b and 2b) (Kerem et al., 1989; Moncalian et al., 2004a) located at the Walker A, D-loop and Rad50 signature motif, respectively (Figure 2a) (Kerem et al., 1989; Moncalian et al., 2004a; De La Rosa Metzere Bierlein & and Scott W. Nelson, 2011).

Mutations at the Walker A domain and Rad50 signature motif (Figure 2a) may also affect important functional sites such as ATP binding site (Table 2). For example, K42R/M/E/A mutation at the Walker A (Figure 1b and 2b) in *S. cerevisiae* and *D. radiodurans* has been identified experimentally to cause defective in ATPase (Table 1) (Chen et al., 2005; Koroleva et al., 2007) and S793R mutation in *Pyrococcus furiosus* showed the inhibition of ATP binding and disrupted communication between ATP loops (Table 1) (Moncalian et al., 2004a). This mutation further distorted the surface of the C-terminal domain and thus altered the interaction between Rad50 monomers to prevent dimerization (Table 1) (Moncalian et al., 2004a). We have also identified mutations at several motifs such as Walker A (G41D, K42M/R/E/A) and Walker B (E1232Q) (Figure 1b, 2a and 2b) that might affect the binding of FHA phosphopeptide ligands

that plays a critical role in DNA damage repair mechanism and cell cycle (Table 2). Many FHA domain-containing proteins localized to the nucleus showed to play a critical role in establishing or maintaining DNA repair, cell cycle checkpoints or transcriptional regulation (Durocher et al., 2000). When mutated, diseases such as Nijmegen breakage syndrome (NBS) and the hereditary cancer syndrome variant Li-Fraumeni (CHK2) will be developed (Matsuura et al., 1998; Varon et al., 1998; Carney et al., 1998b; Featherstone & Jackson, 1998; Bell et al., 1999) suggesting the importance of these conserved residues within Rad50 for DNA repair and maintenance.

Mutations at or near the Rad50 signature motif (Figure 1b and 2a) were also known to be damaging (Table 2), particularly the S1202R (Figure 1b and 2b) mutation which has been studied the most due to its numerous biological defects *in vivo*. The same residue mutations of the Rad50 signature motif in yeast (S12025R) and human (S1202R) also generated complexes that were significantly diminished in adenylate kinase (AK) activity that was important for DNA tethering (Bhaskara et al., 2007). Previously, AK deficiency was found to be associated with anaemia and several cases of mental retardation and psychomotor impairment (Abrusci et al., 2007), which may explain why disruption of the MRN complex also causes this phenotype on patients (Waltes et al., 2009). In addition, such deleterious mutation also contributed to inviable spores and significant telomere shortening in *S. cerevisiae* (Bhaskara et al., 2007). Defects in telomere length in human have been known to cause the pathology of several age-related diseases and premature ageing syndrome, as well as cancer and other human diseases such as Hoyerlaal-Hreidarsson syndrome, Coats plus syndrome, pulmonary fibrosis, dyskeratosis congenita, liver fibrosis and aplastic anemia (Blasco, 2005).

Additionally, most mutations such as G1199E, S1202A/R/M, and Q1205E (Figure 1b and 2b) at the Rad50 signature motif (Figure 2a) were identified to affect PKA phosphorylation site (Table 2) suggesting that this site is dependent upon the function of Rad50 signature motif. Phosphorylation is one of the most ubiquitous and important post translational modifications of proteins, and implicated in almost all kinds of cellular processes and pathways (Ptacek & Snyder, 2006). In neurons, enhanced PKA signalling promotes neuronal development, enhances synaptic plasticity, and elevates dopamine synthesis (Dagda & Das Banerjee, 2015). However a deterioration in PKA signalling has contributed to the aetiology of several neurodegenerative diseases, such as Alzheimer and Parkinson (Dagda & Das Banerjee, 2015). We hypothesized that the defective PKA functional sites may also lead to Nijmegen breakage syndrome associated with neurological phenotype in Rad50 mutations (Waltes et al., 2009), however this potential phosphorylation sites remain to be validated.

C681A and C684R mutations (Figure 1b and 2d) at the zinc hook motif (Figure 2c) were identified deleterious from our analysis (Table 2) and these mutations were known to lead severe defects in various DNA damage response (DDR) such as ataxia-telangiectasia mutated (ATM) protein activation, homologous recombinant, irradiation sensitivity and ataxia telangiectasia and

Rad3 related (ATR) protein activation (He et al., 2012). These findings were consistent with our bioinformatics analysis where C684S deleterious mutation at zinc hook (Figure 1b and 2d) might affect a protein kinase called ataxia telangiectasia mutated (ATM) that belongs to the phosphatidylinositol 3-kinase-related kinase (PIKK) family (Table 2). The ATM protein was known to cause devastating ataxia-telangiectasia syndrome which is characterized by progressive neurological disorder, impaired organ maturation and immunodeficiency (Shiloh & Ziv, 2013). Rad50 phosphorylated ATM at S635 site (Figure 1b and 2d) of which the mutation on this site showed its importance for cell cycle control signalling and DNA repair mechanism (Gatei et al., 2011).

P682E mutation at the zinc hook motif (Figure 1b and 2d) was shown to be deleterious (Table 2), where previous study has reported that the double mutation P682E and S679R at the zinc hook motif have reduced zinc affinity and dimerization efficiency leading to mice lethality (Roset et al., 2014). In addition, crossbreeding P682E and S679R mutant mice with wildtype mice produce offsprings with hydrocephalus (accumulation of cerebrospinal fluid within the brain), defects in hematopoietic stem cells and gametogenic cells. This suggests that the hook motif has strong influence on the MRN complex associated with DDR signalling, tissue homeostasis and tumorigenesis, as well as fertility of the organism (Roset et al., 2014). This is consistent with the mutations in the yeast hook domain that has increased chromosomal fragmentation (Cahill & Carney, 2007), suggesting its presence is required for the binding or tethering of chromosomal ends.

Limitations of *in silico* prediction

Several mutations were functionally predicted to be neutral, in contrast with the previous experimental findings. For example, a few mutations i.e. S635G (*H. sapiens*), S679R, C680N, P682A, V683I (*S. cerevisiae*), V683R (*M. musculus*) and Q685S (*S. cerevisiae*) (Figure 1b and 2d) located at the zinc hook domain (Table 2) and mutations on K22M (*S. cerevisiae* and *M. musculus*), R83I (*S. cerevisiae* and *M. musculus*), T191E, C221E and S106E (*T. maritima*) in the ATPase domain (Figure 1b, 2b and Table 2) were experimentally validated to be deleterious (some causing embryonic lethality, growth defect, cancer predisposition, as well as hematopoietic and spermatogenic depletion *in vivo* (Bender et al., 2002). A few previous studies have also shown discrepancies between computer prediction and experimental data. For example, an extensive *in silico* analysis using PolyPhen2 and MutPred tools of the ATP-binding cassette transporter ABCA1, an important target in anti-atherosclerosis treatment predicted that several nsSNPs can be neutral, contradicting with previous experimental data findings (Marín-Martín et al., 2014). Furthermore, another *in silico* analysis performed using PolyPhen and SIFT on proteins related to several hereditary diseases such as glucose-6-phosphate dehydrogenase deficiency (G6PD), the receptor 1 for tumor necrosis factor-(TNFRSF1A), and familial mediterranean fever (MEFV) has concluded that some nsSNPs impact may also not be predicted

deleterious to correspond to previous phenotypic effect (Tchernitchko, Goossens & Wajcman, 2004). Moreover, *in silico* identification of PmrAB virulence targets in *Salmonella typhimurium* also demonstrated false positive prediction when validated experimentally (Marchal et al., 2004) suggesting that more work has to be done to develop a more accurate bioinformatics prediction platforms in the future. In contrast, various SNP prediction software have predicted that these mutations were not damaging (Table 2).

Such discrepancy between the computational prediction and experimental results may be due to several limitations in the bioinformatic tools used in our analysis. Several web-based prediction tools may supply conflicting results (Wan et al., 2008) and even with an integrated predictor, PredictSNP (Bendl et al., 2014), it is also limited by the differences in algorithms, principles, training datasets and information used. For example, MAPP, PANTHER and SIFT in the PredictSNP used alignment scores for functional prediction whereas SNAP, PoplyPhen-1 and PolyPhen-2 used neural network, support vector machine and Naïve Bayes algorithm, respectively (Bendl et al., 2014). Interestingly, we identified that the software predicts most accurately (in agreement with experimental results) for the motifs or sites located at the highly conserved position (Figure 1b). Conversely, most residues that were predicted to be neutral are located at non-conserved positions in the Rad50 protein (S3 Figure and S7 Table), suggesting that these prediction software may have only been trained and preferentially biased towards conserved regions (Gardner et al., 2017). This suggests that computer prediction should also consider and take into account the effect of non-conserved regions outside the motifs/domains too for future improvement in their algorithms. Furthermore, any subsequent prediction studies should also be aware of this limitation (whether located in conserved or non-conserved regions) to carefully deduce the function of their protein mutation of interest.

Nonetheless, we cannot rule out the possibility that these mutations derived from other organisms may not be readily affected or transferred to other organisms including human. This is because certain organisms may possess gene compensation to compensate or mask the effect of such mutations and that the different proteins from different organisms may not have perfectly superimposable function. Hence, future experiments should focus on their validation especially in human cell line studies to better understand the roles of these mutated residues in Rad50 function.

Conclusions

This study compiled all mutations to date in Rad50 proteins from various organisms and predicts their effects using various software tools such as PredictSNP, MutPred, SNPeffect, I-Mutant and MUpPro. Most predictions for SNPs occurring within conserved regions are in agreement with their corresponding *in vivo* or *in vitro* experimental results. However, SNPs located at non-

conserved regions are less likely to be accurately predicted, and as such algorithms for these software should be improved in future studies. Altogether this study has provided means to prioritized mutations particularly in Rad50 protein that have biologically meaningful function for DNA double-stranded maintenance.

Acknowledgements

The authors would like to acknowledge and thank INBIOSIS and Makmal Genomik 1 research teams from Universiti Kebangsaan Malaysia for their help with laboratory and technical analyses.

References

- Abrusci, P., Chiarelli, L.R., Galizzi, A., Fermo, E., Bianchi, P., Zanella, A. & Valentini, G. (2007). Erythrocyte adenylate kinase deficiency: characterization of recombinant mutant forms and relationship with nonspherocytic hemolytic anemia. *Experimental Hematology*, 35(8): 1182–1189.
- Adzhubei, I., Jordan, D.M. & Sunyaev, S.R. (2013). Predicting functional effect of human missense mutations using PolyPhen-2. *Current Protocols in Human Genetics*, 7(Supp76):20.
- Alani, E., Padmore, R. & Kleckner, N. (1990). Analysis of wild-type and rad50 mutants of yeast suggests an intimate relationship between meiotic chromosome synapsis and recombination. *Cell*, 61(3): 419–436.
- Assenmacher, N. & Hopfner, K.P. (2004). MRE11/RAD50/NBS1: Complex activities. *Chromosoma*, 113(4): 157–166.
- De Baets, G., Van Durme, J., Reumers, J., Maurer-Stroh, S., Vanhee, P., Dopazo, J., Schymkowitz, J. & Rousseau, F. (2012). SNPeffect 4.0: On-line prediction of molecular and structural effects of protein-coding variants. *Nucleic Acids Research*, 40(D1): 935-939.
- Barbi, G., Scheres, J.M.J.C., Schindler, D., Taalman, R.D.F.M., Rodens, K., Mehnert, K., Müller, M., Seyschab, H., Muller, M. & Seyschab, H. (1991). Chromosome instability and X-ray hypersensitivity in a microcephalic and growth-retarded child. *American Journal of Medical Genetics*, 40(1): 44–50.
- Bartoszewski, R. a., Jablonsky, M., Bartoszezwska, S., Stevenson, L., Dai, Q., Kappes, J., Collawn, J.F. & Bebok, Z. (2010). A synonymous single nucleotide polymorphism in delta F508 CFTR alters the secondary structure of the mRNA and the expression of the mutant protein. *Journal of Biological Chemistry*, 285(37): 28741–28748.
- Bell, D.W., Varley, J.M., Szydlo, T.E., Kang, D.H., Wahrer, D.C.R., Shannon, K.E., Lubratovich, M., Verselis, S.J., Isselbacher, K.J., Fraumeni, J.F., Birch, J.M., Li, F.P., Garber, J.E. & Haber, D.A. (1999). Heterozygous germ line hCHK2 mutations in Li-Fraumeni syndrome. *Science*, 286(5449): 2528-2531.

- Bender, C.F., Sikes, M.L., Sullivan, R., Huye, L.E., Le Beau, M.M., Roth, D.B., Mirzoeva, O.K., Oltz, E.M. & Petrini, J.H.J. (2002). Cancer predisposition and hematopoietic failure in Rad50S/S mice. *Genes and Development*, 16(17): 2237–2251.
- Bendl, J., Stourac, J., Salanda, O., Pavelka, A., Wieben, E.D., Zendulka, J., Brezovsky, J. & Damborsky, J. (2014). PredictSNP: Robust and Accurate Consensus Classifier for Prediction of Disease-Related Mutations. *PLoS Computational Biology*, 10(1): 1003440.
- Bhaskara, V., Dupr, A., Lengsfeld, B., Hopkins, B.B., Chan, A., Lee, J.H., Zhang, X., Gautier, J., Zakian, V. & Paull, T.T. (2007). Rad50 Adenylate Kinase Activity Regulates DNA Tethering by Mre11/Rad50 Complexes. *Molecular Cell*, 25(5): 647–661.
- Blasco, M.A. (2005). Telomeres and human disease: Ageing, cancer and beyond. *Nature Reviews Genetics*, 6(8): 611–622.
- van den Bosch, M., Bree, R.T. & Lowndes, N.F. (2003). The MRN complex: Coordinating and mediating the response to broken chromosomes. *EMBO Reports*, 4(9): 844–849.
- Bromberg, Y. & Rost, B. (2007). SNAP: Predict effect of non-synonymous polymorphisms on function. *Nucleic Acids Research*, 35(11): 3823–3835.
- Cahill, D. & Carney, J.P. (2007). Dimerization of the Rad50 protein is independent of the conserved hook domain. *Mutagenesis*, 22(4): 269–274.
- Capriotti, E. & Fariselli, P. (2017). PhD-SNPg: A webserver and lightweight tool for scoring single nucleotide variants. *Nucleic Acids Research*, 45(W1): W247–W252.
- Capriotti, E., Fariselli, P. & Casadio, R. (2005). I-Mutant2.0: Predicting stability changes upon mutation from the protein sequence or structure. *Nucleic Acids Research*, 33(SUPPL. 2): 306–310.
- Cargill, M., Altshuler, D., Ireland, J., Sklar, P., Ardlie, K., Patil, N., Lane, C.R., Lim, E.P., Kalyanaraman, N., Nemesh, J., Ziaugra, L., Friedland, L., Rolfe, A., Warrington, J., Lipshutz, R., Daley, G.Q., Lander, E.S., Cargill, M., Sklar, P., Ireland, J., Ardlie, K., Patil, N., Lane, C.R., Lim, E.P., Kalyanaraman, N., Nemesh, J., Ziaugra, L., Friedland, L., Rolfe, A., Warrington, J., Lipshutz, R. & Daley, G.Q. (1999). Characterization of single-nucleotide polymorphisms in coding regions of human genes. *Nature*, 22(July): 231–238.
- Carney, J.P., Maser, R.S., Olivares, H., Davis, E.M., Le Beau, M., Yates, J.R., Hays, L., Morgan, W.F. & Petrini, J.H. (1998a). The hMre11/hRad50 protein complex and Nijmegen breakage syndrome: linkage of double-strand break repair to the cellular DNA damage response. *Cell*, 93(3): 477–86.
- Carney, J.P., Maser, R.S., Olivares, H., Davis, E.M., Le Beau, M., Yates, J.R., Hays, L., Morgan, W.F. & Petrini, J.H.J. (1998b). The hMre11/hRad50 protein complex and Nijmegen breakage syndrome: Linkage of double-strand break repair to the cellular DNA damage response. *Cell*, 93(3): 477–486.
- Chen, L., Trujillo, K.M., Van Komen, S., Roh, D.H., Krejci, L., Lewis, L.K., Resnick, M.A., Sung, P. & Tomkinson, A.E. (2005). Effect of amino acid substitutions in the rad50 ATP binding domain on DNA double strand break repair in yeast. *The Journal of Biological Chemistry*, 280(4): 2620–2627.

- Cheng, J., Randall, A. & Baldi, P. (2006). Prediction of protein stability changes for single-site mutations using support vector machines. *Proteins*, 62(4): 1125–1132.
- Collins, F.S., Guyer, M.S. & Charkravarti, A. (1997). Variations on a theme: cataloging human DNA sequence variation. *Science*, 278(5343): 1580–1.
- Dagda, R.K. & Das Banerjee, T. (2015). Role of protein kinase A in regulating mitochondrial function and neuronal development: Implications to neurodegenerative diseases. *Reviews in the Neurosciences*, 26(3): 359-370.
- Davidson, A.L., Dassa, E., Orelle, C. & Chen, J. (2008). Structure, function, and evolution of bacterial ATP-binding cassette systems. *Microbiology and Molecular Biology Reviews*, 72(2): 317–364.
- Durocher, D., Taylor, I. a, Sarbassova, D., Haire, L.F., Westcott, S.L., Jackson, S.P., Smerdon, S.J. & Yaffe, M.B. (2000). The molecular basis of FHA domain:phosphopeptide binding specificity and implications for phospho-dependent signaling mechanisms. *Molecular cell*, 6(5): 1169-1182.
- Featherstone, C. & Jackson, S.P. (1998). DNA repair: the Nijmegen breakage syndrome protein. *Current biology*, CB 8(17): R622-5.
- Finn, R.D., Attwood, T.K., Babbitt, P.C., Bateman, A., Bork, P., Bridge, A.J., Chang, H.Y., Dosztanyi, Z., El-Gebali, S., Fraser, M., Gough, J., Haft, D., Holliday, G.L., Huang, H., Huang, X., Letunic, I., Lopez, R., Lu, S., Marchler-Bauer, A., Mi, H., Mistry, J., Natale, D.A., Necci, M., Nuka, G., Orengo, C.A., Park, Y., Pesseat, S., Piovesan, D., Potter, S.C., Rawlings, N.D., Redaschi, N., Richardson, L., Rivoire, C., Sangrador-Vegas, A., Sigrist, C., Sillitoe, I., Smithers, B., Squizzato, S., Sutton, G., Thanki, N., Thomas, P.D., Tosatto, S.C.E., Wu, C.H., Xenarios, I., Yeh, L.S., Young, S.Y. & Mitchell, A.L. (2017). InterPro in 2017-beyond protein family and domain annotations. *Nucleic Acids Research*, 45(D1): D190–D199.
- Gardner, P.P., Paterson, J.M., Ghomi, F.A., Umu, S.U.U., McGimpsey, S. & Pawlik, A. (2017). A meta-analysis of bioinformatics software benchmarks reveals that publication-bias unduly influences software accuracy. *bioRxiv*, 092205.
- Gatei, M., Jakob, B., Chen, P., Kijas, A.W., Becherel, O.J., Gueven, N., Birrell, G., Lee, J.H., Paull, T.T., Lerenthal, Y., Fazry, S., Taucher-Scholz, G., Kalb, R., Schindler, D., Waltes, R., Drk, T. & Lavin, M.F. (2011). ATM protein-dependent phosphorylation of Rad50 protein Regulates DNA repair and cell cycle control. *Journal of Biological Chemistry*, 286(36): 31542–31556.
- Halushka, M.K., Fan, J.B., Bentley, K., Hsie, L., Shen, N., Weder, A., Cooper, R., Lipshutz, R. & Chakravarti, A. (1999). Patterns of single-nucleotide polymorphisms in candidate genes for blood-pressure homeostasis. *Nature Genetics*, 22(3): 239–247.
- He, J., Shi, L.Z., Truong, L.N., Lu, C.S., Razavian, N., Li, Y., Negrete, A., Shiloach, J., Berns, M.W. & Wu, X. (2012). Rad50 zinc hook is important for the Mre11 complex to bind chromosomal DNA double-stranded breaks and initiate various DNA damage responses. *Journal of Biological Chemistry*, 287(38): 31747–31756.

- Hopfner, K.-P., Craig, L., Moncalian, G., Zinkel, R. a, Usui, T., Owen, B. a L., Karcher, A., Henderson, B., Bodmer, J.-L., McMurray, C.T., Carney, J.P., Petrini, J.H.J. & Tainer, J.A. (2002). The Rad50 zinc-hook is a structure joining Mre11 complexes in DNA recombination and repair. *Nature*, 418(6897): 562–566.
- de Jager, M., van Noort, J., van Gent, D.C., Dekker, C., Kanaar, R. & Wyman, C. (2001). Human Rad50/Mre11 is a flexible complex that can tether DNA ends. *Molecular cell*, 8(5): 1129–35.
- Jones, P., O'Mara, M. & George, A. (2009). ABC transporters: a riddle wrapped in a mystery inside an enigma. *Trends in biochemical sciences*, 34(10): 520-31.
- Kelley, L.A., Mezulis, S., Yates, C.M., Wass, M.N. & Sternberg, M.J.E. (2015). The Phyre2 web portal for protein modeling, prediction and analysis. *Nature Protocols*, 10(6): 845–858.
- Kerem, B., Rommens, J.M., Buchanan, J. a, Markiewicz, D., Cox, T.K., Chakravarti, a, Buchwald, M. & Tsui, L.C. (1989). Identification of the cystic fibrosis gene: genetic analysis. *Science*, 245(4922): 1073–1080.
- Koroleva, O., Makharashvili, N., Courcelle, C.T., Courcelle, J. & Korolev, S. (2007). Structural conservation of RecF and Rad50: implications for DNA recognition and RecF function. *The EMBO journal*, 26(3): 867–77.
- De La Rosa Metzere Bierlein & and Scott W. Nelson. (2011). An Interaction between the Walker A and D-loop Motifs Is Critical to ATP Hydrolysis and Cooperativity in Bacteriophage. *The Journal of Biological Chemistry*, 286(29): 26258–26266.
- Lee, J.E., Choi, J.H., Lee, J.H. & Lee, M.G. (2005). Gene SNPs and mutations in clinical genetic testing: Haplotype-based testing and analysis. *Mutation Research - Fundamental and Molecular Mechanisms of Mutagenesis*, 573(1-2): 195-204.
- Luo, G., Yao, M.S., Bender, C.F., Mills, M., Bladl, A.R., Bradley, A. & Petrini, J.H. (1999). Disruption of mRad50 causes embryonic stem cell lethality, abnormal embryonic development, and sensitivity to ionizing radiation. *Proceedings of the National Academy of Sciences of the United States of America*, 96(13): 7376–81.
- Marchal, K., De Keersmaecker, S., Monsieurs, P., van Boxel, N., Lemmens, K., Thijs, G., Vanderleyden, J. & De Moor, B. (2004). In silico identification and experimental validation of PmrAB targets in Salmonella typhimurium by regulatory motif detection. *Genome Biology*, 5(2): R9
- Marín-Martín, F.R., Soler-Rivas, C., Martín-Hernández, R. & Rodríguez-Casado, A. (2014). A Comprehensive In Silico Analysis of the Functional and Structural Impact of Nonsynonymous SNPs in the ABCA1 Transporter Gene. *Cholesterol*, 2014: 1–19.
- Matsuura, S., Tauchi, H., Nakamura, A., Kondo, N., Sakamoto, S., Endo, S., Smeets, D., Solder, B., Belohradsky, B.H., Der Kaloustian, V.M., Oshimura, M., Isomura, M., Nakamura, Y. & Komatsu, K. (1998). Positional cloning of the gene for Nijmegen breakage syndrome. *Nature genetics*, 19(2): 179-181.

- Moncalian, G., Lengsfeld, B., Bhaskara, V., Hopfner, K.P., Karcher, A., Alden, E., Tainer, J.A. & Paull, T.T. (2004). The Rad50 Signature Motif: Essential to ATP Binding and Biological Function. *Journal of Molecular Biology*, 335(4): 937–951.
- Pejaver, V., Urresti, J., Lugo-Martinez, J., Pagel, K.A., Lin, G.N., Nam, H.-J., Mort, M., Cooper, D.N., Sebat, J., Iakoucheva, L.M., Mooney, S.D. & Radivojac, P. (2017). MutPred2: inferring the molecular and phenotypic impact of amino acid variants. *bioRxiv*, :134981.
- Pettersen, E.F., Goddard, T.D., Huang, C.C., Couch, G.S., Greenblatt, D.M., Meng, E.C. & Ferrin, T.E. (2004). UCSF Chimera--a visualization system for exploratory research and analysis. *Journal Comput Chemistry*, 25(13): 1605–1612.
- Piro, R.M. & Di Cunto, F. (2012). Computational approaches to disease-gene prediction: Rationale, classification and successes. *FEBS Journal*, 279(5): 678-696.
- Ptacek, J. & Snyder, M. (2006). Charging it up: global analysis of protein phosphorylation.
- Ramensky, V. 2002. Human non-synonymous SNPs: server and survey. *Nucleic Acids Research*, 30(17): 3894–3900.
- Risch, N. & Merikangas, K. (1996). The future of genetic studies of complex human diseases. *Science*, 273(5281): 1516–1517.
- Rojowska, A., Lammens, K., Seifert, F.U., Drenth, C. & Feldmann, H. (2014). Structure of the Rad 50 DNA double-strand break repair protein in complex with DNA. *The EMBO journal*, 33(23): 2847–2859.
- Roset, R., Inagaki, A., Hohl, M., Brenet, F., Lafrance-Vanasse, J., Lange, J., Scandura, J.M., Tainer, J.A., Keeney, S. & Petrini, J.H.J. (2014). The Rad50 hook domain regulates DNA damage signaling and tumorigenesis. *Genes and Development*, 28(5): 451–462.
- Schymkowitz, J., Borg, J., Stricher, F., Nys, R., Rousseau, F. & Serrano, L. (2005). The FoldX web server: An online force field. *Nucleic Acids Research*, 33(SUPPL. 2).
- Shaikho, E.M., Farrell, J.J., Alsultan, A., Qutub, H., Al-Ali, A.K., Figueiredo, M.S., Chui, D.H.K., Farrer, L.A., Murphy, G.J., Mostoslavsky, G., Sebastiani, P. & Steinberg, M.H. (2017). A phased SNP-based classification of sickle cell anemia HBB haplotypes. *BMC Genomics*, 18(1): 608.
- Shiloh, Y. & Ziv, Y. (2013). The ATM protein kinase: Regulating the cellular response to genotoxic stress, and more. *Nature reviews. Molecular cell biology*, 14(4): 197-210.
- Sievers, F. & Higgins, D.G. (2018). Clustal Omega for making accurate alignments of many protein sequences. *Protein Science*, 27(1):135-145.
- Sim, N.-L., Kumar, P., Hu, J., Henikoff, S., Schneider, G. & Ng, P.C. (2012). SIFT web server: predicting effects of amino acid substitutions on proteins. *Nucleic Acids Research*, 40(W1): W452–W457.
- Stone, E.A. & Sidow, A. (2005). Physicochemical constraint violation by missense substitutions mediates impairment of protein function and disease severity. *Genome Research*, 15(7): 978–986.
- Symington, L.S. (2002). Role of RAD52 epistasis group genes in homologous recombination and double-strand break repair. *Microbiology and molecular biology reviews*, 66(4): 630–70.

792 Tchernitchko, D., Goossens, M. & Wajcman, H. (2004). In silico prediction of the deleterious
793 effect of a mutation: Proceed with caution in clinical genetics. *Clinical Chemistry*, 50(11):
794 1974-1978.

795 Thompson, J.D., Gibson, T.J. & Higgins, D.G. (2002). Multiple Sequence Alignment Using
796 ClustalW and ClustalX. *Current Protocols in Bioinformatics*, 0(1): 2.3.1-2.3.22.

797 Traeger, J., Wood, W.G., Clegg, J.B., Weatherall, D.J. & Wasi, P. (1980). Defective synthesis of
798 HbE is due to reduced levels of β E mRNA. *Nature* 288(5790): 497-499.

799 Tranchevent, L.C., Capdevila, F.B., Nitsch, D., de Moor, B., de Causmaecker, P. & Moreau, Y.
800 (2011). A guide to web tools to prioritize candidate genes. *Briefings in Bioinformatics*,
801 12(1): 22-32.

802 Varon, R., Vissinga, C., Platzer, M., Cerosaletti, K.M., Chrzanowska, K.H., Saar, K., Beckmann,
803 G., Seemanová, E., Cooper, P.R., Nowak, N.J., Stumm, M., Weemaes, C.M., Gatti, R. a,
804 Wilson, R.K., Digweed, M., Rosenthal, a, Sperling, K., Concannon, P. & Reis, a. (1998).
805 Nibrin, a novel DNA double-strand break repair protein, is mutated in Nijmegen breakage
806 syndrome. *Cell*, 93(3): 467-476.

807 Waltes, R., Kalb, R., Gatei, M., Kijas, A.W., Stumm, M., Sobeck, A., Wieland, B., Varon, R.,
808 Lerenthal, Y., Lavin, M.F., Schindler, D. & Dörk, T. (2009). Human RAD50 Deficiency in a
809 Nijmegen Breakage Syndrome-like Disorder. *American Journal of Human Genetics*, 84(5):
810 605-616.

811 Wan, J., Kang, S., Tang, C., Yan, J., Ren, Y., Liu, J., Gao, X., Banerjee, A., Ellis, L.B.M. & Li,
812 T. (2008). Meta-prediction of phosphorylation sites with weighted voting and restricted grid
813 search parameter selection. *Nucleic Acids Research*, 36(4): 22.

Figure 1

Domain analysis and multiple sequence alignment.

Domain analysis using InterPro shows that Rad50 contains P-loop containing nucleoside triphosphate hydrolase domain belongs to the ATP Binding Cassette (ABC) protein superfamily (pink box) as well as a special domain called zinc hook, which so far does not overlap with any homologous superfamilies (blue box) (a). ABC protein consists of six conserved motifs; i.e. Walker A (WA), Q-loop (QL), signature motif (SM), Walker B (WB), D-loop (DL), and H-loop (HL) which make up the nucleotide binding domain. Zinc hook domain contains a conserved CxxC motif located at the residue number 681-684 (a). All deleterious residues identified from the literature were highlighted based on human equivalent mutation (Supplementary Table S1) and those occurring only in the conserved regions are shown in (b). Multiple sequence alignment (MSA) analysis of Rad50 sequences dataset (human, *D. rerio* (zebrafish), *M. musculus* (mouse), *R. norvegicus* (rat), *P. troglodytes* (chimpanzee), *O. cuniculus* (rabbit), *D. melanogaster* (fruit fly), *S. cerevisiae* (yeast) and *S. pombe* (yeast)) showed conserved residues in specific motifs (b). An “*” (asterisk) indicates position which has a single, fully conserved residue. A “:” (colon) indicates conservation between groups of strongly similar properties - scoring > 0.5 in the Gonnet PAM250 matrix. A “.” (period) indicates conservation between groups of weakly similar properties - scoring ≤ 0.5 in the Gonnet PAM250 matrix.

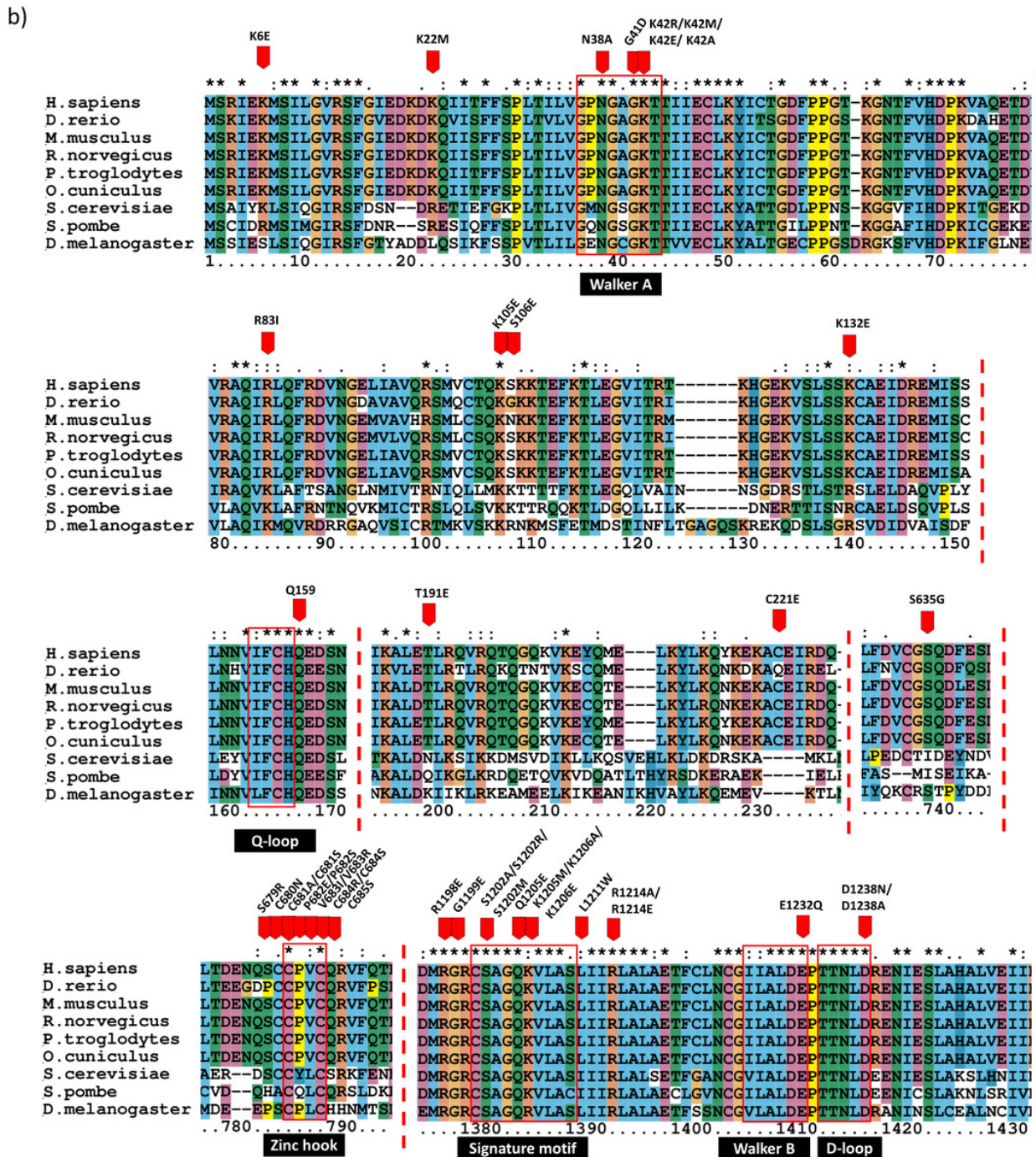
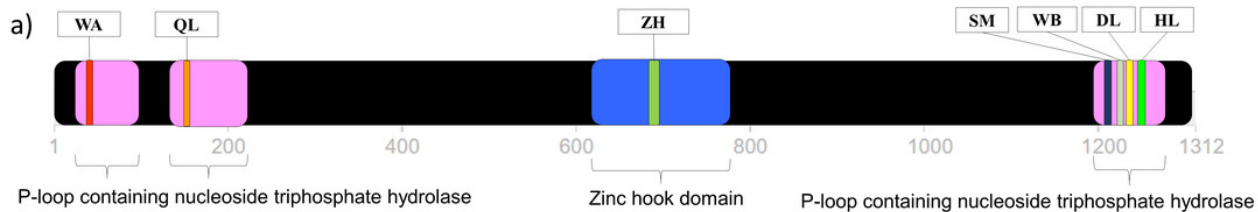


Figure 2

3D structure of Rad50.

3D structure of Rad50. A 3D structure of Rad50 human modelled using fold recognition technique Phyre2 using structure from *Chaetomium thermophilum* as a template (PDB ID: 5DAC). N-terminal of 276 residues (2-278) and C-terminal of 155 residues (1153-1306) are colored as blue and green; respectively (a). All six motifs identified are marked and represented by ball and stick representation with different colours (orange for Q-loop, blue for Walker B, purple for signature motif, yellow for D-loop, grey for Walker A and green from H-loop) (a). All mutated residues identified were marked and labelled in the 3D structure (b). Zinc hook structure of 181 residues (residue 585-766) that has been determined (PDB ID: 5GOX) and its deleterious residues also marked in the structure (c,d). All figures were generated using UCSF Chimera (Pettersen et al., 2004)

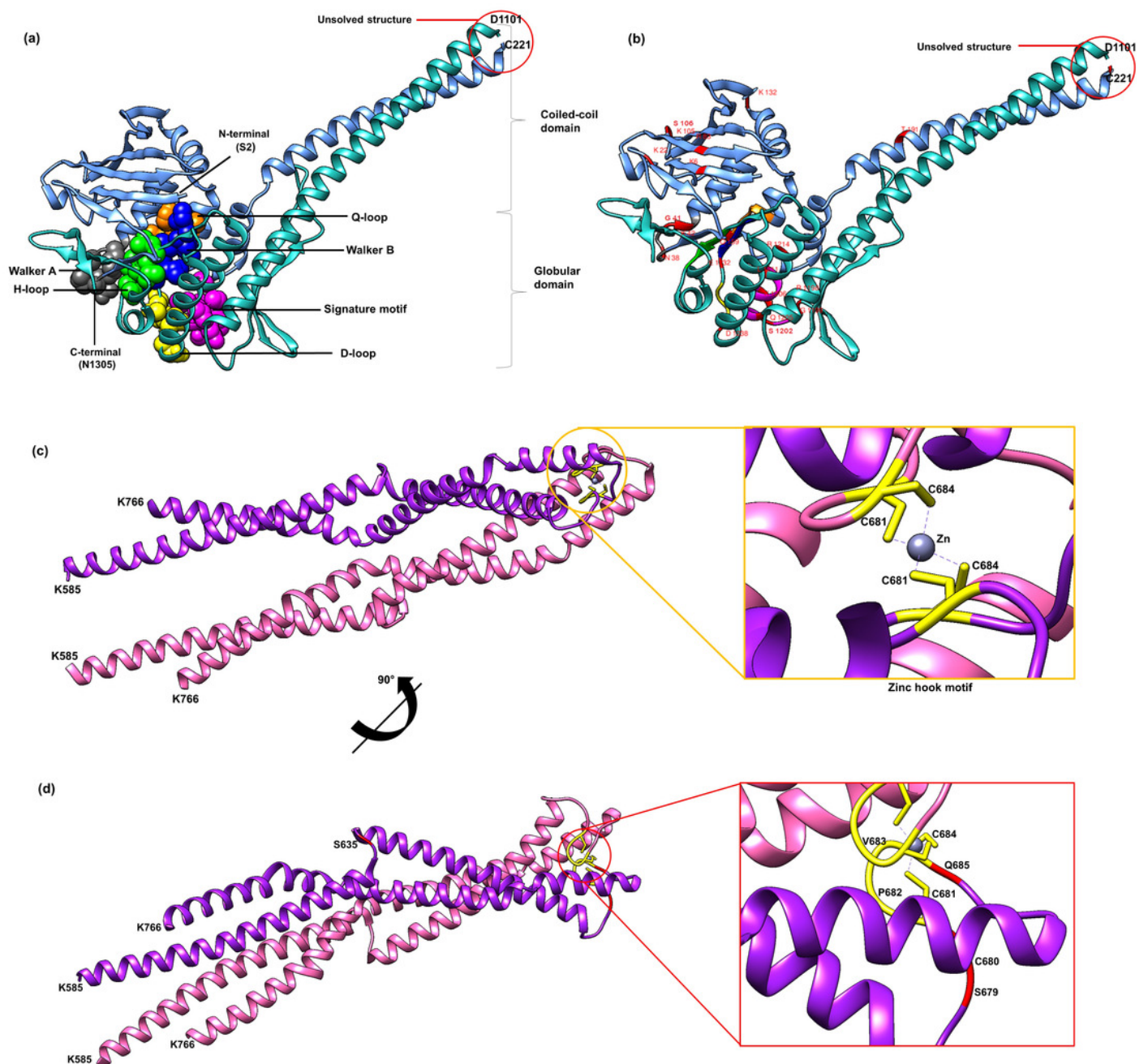


Table 1(on next page)

Summary of the most damaging effects of Rad50 mutations obtained from previous *in vitro* and *in vivo* experiments.

Abbreviations: HR (homologous recombination), NHEJ (non-homologous end joining repair), CFTR (cystic fibrosis transmembrane conductance regulator), ATP (adenosine tri-phosphate), ATM (ataxia-telangiectasia mutated), ATR (ATM-and Rad3-Related). Refer to Table S2 for the description of all mutations.

1

Motif/domain	Mutations	Organism	Effects	Ref.
Walker A	K40A/R/E	<i>S. cerevisiae</i>	<ul style="list-style-type: none"> • HR and NHEJ defects and lower ATPase activity 	(Chen <i>et al.</i> , 2005)
Walker A D-loop	N38A, D512N/A	T4 bacteriophage	<ul style="list-style-type: none"> • Naturally occurring mutation of CFTR protein • Reduce in ATP activity 	(De La Rosa Metzere Bierlein and and Scott W. Nelson, 2011)
ATP binding domain and Walker A	G39D, K40E K81I, R20M	<i>S. cerevisiae</i>	<ul style="list-style-type: none"> • Total defect in formation of viable spore 	(Alani, Padmore and Kleckner, 1990)
ATP binding domain	K6E, K22M, R83I	<i>M. musculus</i>	<ul style="list-style-type: none"> • Embryonic lethality, growth defect, cancer predisposition, hematopoietic and spermatogenic depletion 	(Bender <i>et al.</i> , 2002)
Walker A	K39R, K42M	<i>D. radiodurans</i>	<ul style="list-style-type: none"> • Prevented ATP binding and hydrolysis 	(Koroleva <i>et al.</i> , 2007)
ATPase binding domain, Walker B and Signature motif	K115E, K175E, K182E, R94E, K95E, R765E E798Q, S768R, K103E, K104E, R131E, R1202E, S1205R, E1235Q	<i>T. maritima</i> <i>S. cerevisiae</i>	<i>In vitro: Thermotoga maritima</i> <ul style="list-style-type: none"> • K175E, K182E, K115E Reduced DNA binding • R94E and K95E: Important for DNA binding • R765E: Diminished DNA binding • E798Q: Low affinity to DNA • S768R: Reduced DNA binding <i>In vivo: Saccharomyces cerevisiae</i> <ul style="list-style-type: none"> • S1205R and E1235Q double mutation: Unable to rescue the impaired DNA damage response • K103E, K104E and R131E: Strongly affected DNA binding and moderate reduction in telomere length • K103E and R131R (double mutation) and R1201E: Significantly reduced telomere length • S1205R: Significantly reduced telomere length 	(Rojowska <i>et al.</i> , 2014b)
Zinc hook	S679R, P682E V683R	<i>M. musculus</i>	<ul style="list-style-type: none"> • Lethality in mice. Hydrocephalus, defects in primitive hematopoietic and gametogenic cells 	(Roset <i>et al.</i> , 2014)
	C684N, C685A P686A V687I C688R Q689S	<i>S. cerevisiae</i>	<ul style="list-style-type: none"> • Defective to be recruited to chromosomal double strand break • Phenotype as severe as Rad50 null mutant • Defective in ATM activation, HR, sensitive to irradiation and ATR activation 	(He <i>et al.</i> , 2012)
	C288S, C291S	T4 bacteriophage	<ul style="list-style-type: none"> • Double mutation is lethal 	(Barfoot <i>et al.</i> , 2015)
	S635G	<i>H. sapiens</i>	<ul style="list-style-type: none"> • Chromosomal instability • Defective ATM-dependent signalling 	(Gatei <i>et al.</i> , 2011)

	S685R Y688E L689R	<i>S. cerevisiae</i>	<ul style="list-style-type: none"> • S685R and Y688E double mutation: Sporulation efficiency and viability were severely impaired followed by L689R • Rad50-Mre11 interaction was strongly impaired, partial suppression of telomere and meiotic defects 	(Hohl <i>et al.</i> , 2015)
Rad50 Signature motif	R805E L802W	<i>P. furiosus</i>	<ul style="list-style-type: none"> • L802W: Decrease dimerization in ATP, hydrolysis and cleavage site • R805E: Poorly grown in camptothecin; inability to repair endogenous DNA damage by HR and showed defect in resection in HO endonuclease induced 	(Deshpande <i>et al.</i> , 2014)
	K1187A K1187E R1195A R1195E	<i>S. pombe</i>	<ul style="list-style-type: none"> • K1187A: Sensitive in higher dose of clastogens • K1187E, R1195A and R1195E: Significantly sensitive to clastogen agents and were deleterious as Rad50 null mutation 	(Williams <i>et al.</i> , 2011)
	S471A/R/M, E474Q,K475M	T4 bacteriophage	<ul style="list-style-type: none"> • S471A/R.M, E474Q and K475M: Residues involved in the allosteric transmission between DNA and ATP binding sites 	(Herdendorf and Nelson, 2011)
	S1205R S793R S1202R	<i>S. cerevisiae</i> <i>P. furiosus</i> <i>H. sapiens</i>	<ul style="list-style-type: none"> • S1202R: Reduced adenylate kinase • S793R: Deficient in ATP-dependent dimer formation and ATP binding • S1202R and S1205R: Low level of adenylate kinase • S1205R: Telomere shortening, not support spore viability 	(Bhaskara <i>et al.</i> , 2007)
Signature motif and Q loop	S793R Q140H	<i>P. furiosus</i>	<ul style="list-style-type: none"> • S793R: Analogs to the mutation in CFTR (S549R) gene that results cystic fibrosis • S793R: Prevented ATP binding 	(Moncalian, Lengsfeld, Bhaskara, Hopfner, Karcher, Alden, J. A. Tainer, <i>et al.</i> , 2004)
	S1205R	<i>S. cerevisiae</i>	<ul style="list-style-type: none"> • S1205R <i>S. cerevisiae</i>: Failed to complement Rad50 deletion strain in DNA repair assay • S783R and Q140H: Halted ATP-dependent activities 	
ATPase domain	R1093(stop) c.3939A/T	<i>H. sapiens</i>	<ul style="list-style-type: none"> • Nijmegen breakage syndrome like disorder (NBSLD) 	(Waltes <i>et al.</i> , 2009)

Table 2 (on next page)

In silico analysis of 42 deleterious mutations in Rad50.

Different tools were used to analyse all mutations as abbreviated in the table. PS (PredictSNP), MP (MAPP), PhS (PhD-SNP), PP1 (Poly-Phen1), PP2 (Poly-Phen2), SF (SIFT), SN (SNAP), IM (I-Mutant), MPr (MuPro), Pr (probability), AG (protein aggregation) and AM (Amyloid aggregation). Please refer to Materials and Methods for detailed descriptions of these tools. Note that all mutations listed above are based on the equivalent mutations in human.

1

		Amino acid Impact (PredictSNP) (Neutral/Deleterious)							Molecular mechanisms (MutPred2)				Structural phenotyping (SNPeffect)		Protein stability (Imutant/MuPro)	
Motif	Mutation [source]	PS	MP	Ph-S	PP-1	PP-2	SF	SP	Affected molecular mechanisms	Pr	P-value	Affected functional sites	AG	AM	IM	MPr
Walker A	N38A (De La Rosa Metzere Bierlein and Scott W. Nelson, 2011)	D	D	D	D	D	D	D	Loss of catalytic site at N38 Loss of relative solvent accessibility Altered ordered interface Loss of allosteric site at N38 Altered DNA binding Altered metal binding Gain of methylation at K42	0.53 0.34 0.33 0.31 0.25 0.23 0.17	8.8e-05 3.4e-03 9.9e-03 3.0e-03 6.7e-03 0.02 8.4e-03	• ATP/GTP-binding site motif A (P-loop)	No	No	↓	↓
	G41D (Alani, Padmore and Kleckner, 1990)	D	D	D	D	D	D	D	Altered metal binding Gain of allosteric site at G41 Altered ordered interface Gain of relative solvent accessibility Gain of helix Loss of strand Altered DNA binding Loss of catalytic site at N38 Loss of methylation at K42	0.40 0.33 0.30 0.29 0.29 0.28 0.28 0.27 0.17	1.7e-04 6.4e-04 0.02 0.01 0.01 9.7e-03 4.3e-03 3.7e-03 0.01	• FHA phosphopeptide ligands • CK2 Phosphorylation site • N-myristoylation site • ATP/GTP-binding site motif A (P-loop)	No	No	↓	↓
	K42R (Chen <i>et al.</i> , 2005) (Koroleva <i>et al.</i> , 2007)	D	D	D	D	D	D	D	Loss of relative solvent accessibility Altered DNA binding Loss of allosteric site at T44 Loss of catalytic site at N38 Altered metal binding Loss of methylation at K42	0.30 0.28 0.28 0.27 0.25 0.20	9.6e-03 3.2e-03 7.3e-03 3.4e-03 0.01 6.9e-03	• FHA phosphopeptide ligands • PKA phosphorylation site • CK2 phosphorylation site • N-myristoylation site • ATP/GTP-binding site motif A (P-loop)	No	No	↓	↓
	K42M (Koroleva <i>et al.</i> , 2007)	D	D	D	D	D	D	D	Altered DNA binding Loss of allosteric site at K42 Loss of relative solvent accessibility Altered ordered interface Gain of catalytic site at T43 Altered metal binding Loss of methylation at K42	0.37 0.36 0.34 0.32 0.29 0.28 0.20	7.1e-04 1.1e-03 3.1e-03 0.01 1.6e-03 5.8e-03 6.8e-03	• FHA phosphopeptide ligands • CK2 phosphorylation site • N-myristoylation site • ATP/GTP-binding site motif A (P-loop)	No	No	↑	↑
	K42E (Alani, Padmore and Kleckner, 1990) (Chen <i>et al.</i> , 2005)	D	D	D	D	D	D	D	Altered metal binding Gain of catalytic site at T43 Altered DNA binding Loss of allosteric site at K42 Altered ordered interface Loss of relative solvent accessibility Gain of strand Loss of methylation at K42	0.44 0.33 0.33 0.30 0.29 0.29 0.27 0.20	4.5e-04 6.4e-04 1.3e-03 4.3e-03 0.02 0.01 0.03 6.8e-03	• FHA phosphopeptide ligands • CK2 phosphorylation site • Polo-like kinase phosphorylation sit • N-myristoylation site • ATP/GTP-binding site motif A (P-loop)	No	No	↓	↓
	K42A (Chen <i>et al.</i> , 2005)	D	D	D	D	D	D	D	Loss of allosteric site at K42 Loss of relative solvent accessibility Altered DNA binding Altered ordered interface Gain of catalytic site at T43 Altered metal binding Loss of methylation at K42	0.53 0.37 0.37 0.36 0.32 0.31 0.20	5.4e-05 1.6e-03 6.0e-04 4.2e-03 8.7e-04 0.01 6.8e-03	• FHA phosphopeptide ligands • CK2 phosphorylation site • N-myristoylation site • ATP/GTP-binding site motif A (P-loop)	No	No	↓	↓
	Q-loop	D	D	D	D	D	D	D	No effect	-	-	None	No	No	↓	↓

Zinc hook	<i>al.</i> , 2004)																
	S635G (Gatei <i>et al.</i> , 2011)	N	N	N	D	N	N	N	No effect	-	-	None	No	No	↓	↓	
	S679R (Roset <i>et al.</i> , 2014)(Hohl <i>et al.</i> , 2015)	N	N	N	N	N	D	N	No effect	-	-	None	No	No	↓	↓	
	C680N (He <i>et al.</i> , 2012)	N	N	N	D	N	D	D	Loss of N-linked glycosylation at N677	0.02	0.04	• N-glycosylation site	No	No	↓	↓	
	C681A (He <i>et al.</i> , 2012)	D	D	D	D	D	D	D	Gain of helix Gain of N-linked glycosylation at N677	0.30 0.02	8.0e-03 0.03	• N-glycosylation site	No	No	↓	↓	
	C681S (Barfoot <i>et al.</i> , 2015)	D	D	D	D	D	D	N	Gain of N-linked glycosylation at N677	0.02	0.04	• N-glycosylation site • Proline-directed phosphorylation • MAPK phosphorylation site	No	No	↑	↓	
	P682E (Roset <i>et al.</i> , 2014) (Hohl <i>et al.</i> , 2015)	D	D	N	D	D	D	D	Gain of helix Altered coiled coil Loss of N-linked glycosylation at N677	0.32 0.14 0.02	3.1e-03 0.03 0.04	• CK2 phosphorylation site	No	No	↓	↓	
	P682A (He <i>et al.</i> , 2012)	N	N	N	N	D	D	N	No effect	-	-	None	No	No	↓	↓	
	V683I (He <i>et al.</i> , 2012)	N	N	N	N	N	D	N	No effect	-	-	None	No	No	↑	↑	
	V683R (Roset <i>et al.</i> , 2014)(Hohl <i>et al.</i> , 2015)	N	N	N	N	N	D	D	Gain of helix	0.29	0.01	None	No	No	↓	↓	
	C684R (He <i>et al.</i> , 2012)	D	D	D	D	D	D	D	Gain of helix	0.30	9.5e-03	• Diarginine retention/retrieving signal	No	No	↓	↓	
	C684S (Barfoot <i>et al.</i> , 2015)	D	D	D	D	D	D	D	No effect	-	-	• PIKK phosphorylation site • PKC phosphorylation site	No	No	↓	↓	
Q685S (He <i>et al.</i> , 2012)	N	N	N	N	N	D	N	Altered coiled coil	0.53	6.5e-03	• BRCT phosphopeptide ligands • USP7 binding motif	No	No	↓	↓		
Signature motif	R1198E (Rojowska <i>et al.</i> , 2014b)	D	D	D	D	D	D	D	Gain of catalytic site at R1200 Gain of allosteric site at R1200 Altered metal binding Altered transmembrane protein	0.25 0.21 0.14 0.10	4.8e-03 0.03 0.03 0.04	• Diarginine retention/retrieving signal	No	No	↓	↓	
	G1199E (Rojowska <i>et al.</i> , 2014b)	D	D	D	D	D	D	D	Loss of allosteric site at R1200 Loss of catalytic site at R1200 Altered transmembrane protein	0.23 0.20 0.11	0.03 0.01 0.04	• Diarginine retention/retrieving signal • PKA phosphorylation site	No	No	↓	↓	
	S1202A (Herdendorf and Nelson, 2011)	D	D	D	D	D	D	D	Loss of allosteric site at R1200 Loss of catalytic site at R1200	0.23 0.20	0.03 0.01	• PKA phosphorylation site • -Glycosaminoglycan attachment site	No	No	↓	↓	
	S1202R (Koroleva <i>et al.</i> , 2007) (Moncalian, Lengsfeld, Bhaskara, Hopfner, Karcher, Alden, J. A. Tainer, <i>et al.</i> , 2004) (Herdendorf and Nelson, 2011)	D	D	D	D	D	D	D	Gain of ADP-ribosylation at S1202 Loss of allosteric site at R1200 Loss of catalytic site at R1200	0.25 0.23 0.21	8.4e-03 0.03 1.0e-02	• PKA phosphorylation site • Glycosaminoglycan attachment site	No	No	↑	↓	

	(Bhaskara <i>et al.</i> , 2007)															
	S1202M (Herdendorf and Nelson, 2011)	D	D	D	D	D	D	D	Loss of allosteric site at R1200 Gain of catalytic site at R1200	0.23 0.21	0.02 9.7e-03	<ul style="list-style-type: none"> PKA phosphorylation site Glycosaminoglycan attachment site 	No	No	↑	↑
	Q1205E (Herdendorf and Nelson, 2011)	D	D	D	N	D	D	D	Gain of allosteric site at R1200 Loss of catalytic site at R1200	0.23 0.20	0.02 0.01	<ul style="list-style-type: none"> PKA phosphorylation site CK2 phosphorylation site 	↑	↓	↓	↓
	K1206M (Herdendorf and Nelson, 2011)	D	D	D	D	D	D	D	Gain of catalytic site at S1202	0.09	0.04	None	↑	↓	↑	↑
	K1206A (Williams <i>et al.</i> , 2011)	D	D	D	D	D	D	D	Loss of catalytic site at K1206	0.09	0.05	None	↑	↓	↓	↑
	K1206E (Williams <i>et al.</i> , 2011)	D	D	D	D	D	D	D	Gain of catalytic site at K1206	0.11	0.03	<ul style="list-style-type: none"> TRAF2 binding site NES nuclear export signal 	↑	↓	↓	↑
	L1211W (Deshpande <i>et al.</i> , 2014)	D	D	D	D	D	D	D	Loss of catalytic site at K1206	0.08	0.05	<ul style="list-style-type: none"> SUMO interaction site 	↑	↑	↓	↓
	R1214A (Williams <i>et al.</i> , 2011)	D	D	D	D	D	D	D	Loss of allosteric site at R1214	0.22	0.03	<ul style="list-style-type: none"> ATP-binding cassette, ABC transporter-type, signature and profile 	↑	↓	↓	↓
	R1214E (Deshpande <i>et al.</i> , 2014) (Williams <i>et al.</i> , 2011)	D	D	D	D	D	D	D	Loss of allosteric site at R1214	0.20	0.04	<ul style="list-style-type: none"> SUMO interaction site ATP-binding cassette, ABC transporter-type, signature and profile 	↑	No	↓	↓
Walker B	E1232Q (Rojowska <i>et al.</i> , 2014b)	D	D	D	D	D	D	D	Altered metal binding Loss of catalytic site at E1232 Loss of allosteric site at P1233 Altered transmembrane protein	0.48 0.34 0.24 0.12	4.3e-03 9.8e-04 0.02 0.03	<ul style="list-style-type: none"> FHA phosphopeptide ligands SUMO interaction site 	No	No	↓	↓
D-loop	D1238N (De La Rosa Metzere Bierlein and and Scott W. Nelson, 2011)	D	D	D	D	D	D	D	Altered ordered interface Altered metal binding Gain of relative solvent accessibility Gain of allosteric site at P1233 Loss of catalytic site at T1234 Altered transmembrane protein Altered coiled coil	0.3 0.31 0.27 0.25 0.17 0.12 0.08	4.3e-03 2.8e-03 0.02 0.01 0.02 0.03 0.05	<ul style="list-style-type: none"> FHA phosphopeptide ligands Casein kinase II phosphorylation site 	No	No	↓	↓
	D1238A (De La Rosa Metzere Bierlein and and Scott W. Nelson, 2011)	D	D	D	D	D	D	D	Altered metal binding Altered ordered interface Loss of allosteric site at P1233 Loss of catalytic site at T1234 Altered transmembrane protein	0.41 0.40 0.26 0.18 0.12	3.4e-04 1.4e-03 0.01 0.02 0.02	<ul style="list-style-type: none"> FHA phosphopeptide ligands Casein kinase II phosphorylation site 	No	No	↓	↓
ATPase domain/ coiled-coil	K6E (Alani, Padmore and Kleckner, 1990) (Bender <i>et al.</i> , 2002)	D	D	D	D	D	N	D	Loss of strand Altered DNA binding Gain of N-terminal acetylation at M1	0.27 0.16 0.03	0.03 0.04 4.1e-03	None	No	No	↓	↓
	K22M (Alani, Padmore and Kleckner, 1990) (Bender	N	N	N	N	N	D	N	No effect	-	-	None	↑	No	↓	↑

	<i>et al.</i> , 2002)															
	R83I (Alani, Padmore and Kleckner, 1990) (Bender <i>et al.</i> , 2002)	N	D	N	N	N	D	N	Altered ordered interface Altered DNA binding Altered coiled coil	0.29 0.22 0.10	0.03 0.02 0.04	• PP1-docking motif RVXF	↑	↑	↓	↑
	K132E (Rojowska <i>et al.</i> , 2014b)	D	D	N	D	D	D	D	Loss of helix Altered transmembrane protein Gain of strand	0.28 0.27 0.27	0.02 7.3e-04 0.01	• CK1 phosphorylation site • Protein kinase C phosphorylation site	No	No	↓	↓
	T191E (Rojowska <i>et al.</i> , 2014b)	N	D	N	N	N	N	N	Altered coiled coil Loss of acetylation at K187	0.28 0.28	0.01 6.2e-03	• TRAF2 binding site • NEK2 phosphorylation site • PKC phosphorylation site	No	No	↓	↓
	C221E (Rojowska <i>et al.</i> , 2014b)	N	N	N	N	N	N	N	No effect	-	-	None	No	No	↓	↓
	K105E (Rojowska <i>et al.</i> , 2014b)	D	N	D	D	D	D	D	No effect	-	-	None	No	No	↓	↓
	S106E (Rojowska <i>et al.</i> , 2014b)	N	N	N	N	N	N	N	No effect	-	-	None	No	No	↑	↓

