

Minimum Information Specification For In Situ Hybridization and Immunohistochemistry Experiments (MISFISHIE)

Eric W Deutsch¹, Catherine A Ball², Jules J Berman³, G Steven Bova⁴, Alvis Brazma⁵, Roger E Bumgarner⁶, David Campbell¹, Helen C Causton⁷, Jeff Christiansen⁸, Duncan R Davidson⁸, Young Ah Goo^{1,10}, Sean Grimmond¹¹, Thorsten Henrich¹², Bernhard G Herrmann¹³, Michael H Johnson¹, Martin Korb¹, Jason C Mills¹⁴, Asa J Oudes^{1,10}, Helen E Parkinson⁵, Laura E Pascal^{1,10}, Nicolas Pollet¹⁵, John Quackenbush¹⁶, Mirana Ramialison¹², Martin Ringwald¹⁷, Susanna-A Sansone⁵, Gavin Sherlock¹⁸, Christian J Stoeckert, Jr.¹⁹, Jason Swedlow²⁰, Ronald C Taylor²¹, Laura Walashek^{1,10}, Anthony Warford²², David G Wilkinson²³, Yi Zhou²⁴, Leonard I Zon²⁵, Alvin Y Liu^{1,10}, Lawrence D True²⁶

¹ Institute for Systems Biology, 1441 N 34th Street, Seattle, WA 98103, USA

² Department of Biochemistry, Stanford University School of Medicine, 279 Campus Drive West, Stanford, CA 94305, USA

³ Association for Pathology Informatics, 9650 Rockville Pike, Bethesda, MD 20814, USA

⁴ Johns Hopkins University School of Medicine, Departments of Pathology, Health Information Sciences, Genetic Medicine, Oncology, and Urology, Baltimore, MD, USA

⁵ EMBL-EBI, Wellcome Trust Genome Campus, Hinxton, Cambridge CB10 1SD, UK

⁶ Department of Microbiology, University of Washington, Seattle, WA 98195, USA

⁷ CSC/Imperial College School of Medicine Microarray Centre, London, W12 ONN, UK

⁸ MRC Human Genetics Unit, Western General Hospital, Crewe Road, Edinburgh, UK EH4 2XU

¹⁰ Department of Urology, University of Washington, Seattle, WA 98195, USA

¹¹ Institute for Molecular Bioscience, University of Queensland, St. Lucia, Qld, Australia

¹² EMBL Heidelberg, Meyerhofstrasse 1, D-69117 Heidelberg, Germany

¹³ Institute for medical Genetics, Charité-CBF, and Dept. Developmental Genetics, Max-Planck-Institute for molecular Genetics, Ihnestrasse 73, 14195 Berlin, Germany

¹⁴ Department of Pathology and Immunology, Washington University School of Medicine, St. Louis, MO 63110, USA

¹⁵ CNRS UMR 8080, Université Paris-Sud, 91405 Orsay, France

¹⁶ Dana-Farber Cancer Institute, 44 Binney Street, M232, Boston, MA 02115, USA

¹⁷ The Jackson Laboratory, 600 Main Street, Bar Harbor, ME 04609, USA

¹⁸ Department of Genetics, Stanford University School of Medicine, 300 Pasteur Drive, Stanford, CA 94305-5120, USA

¹⁹ Center for Bioinformatics and Department of Genetics, University of Pennsylvania, Philadelphia, PA 19104, USA

²⁰ Division of Gene Regulation and Expression, Wellcome Trust Biocentre, University of Dundee, Dow Street, Dundee DD1 5EH, Scotland

²¹ Pacific Northwest National Laboratory, PO Box 999, MS K7-90, Richland, WA 99352, USA

²² Wellcome Trust Sanger Institute, Hinxton, Cambridge CB10 1HH, UK

²³ Division of Developmental Neurobiology, National Institute for Medical Research, The Ridgeway, Mill Hill, London NW7 1AA, UK

²⁴ Children's Hospital Boston and Harvard Medical School, Boston, MA 02115, USA

²⁵ Division of Hematology/Oncology, Children's Hospital, Karp Research Laboratories, 1 Blackfan Circle, Boston, Massachusetts 02115, USA

²⁶ Department of Pathology, University of Washington, Seattle, WA 98195-6100, USA

Abstract

One purpose of the biomedical literature is to report results in sufficient detail so that the methods of data collection and analysis can be independently replicated and verified. Here we present for consideration a minimum information specification for gene expression localization experiments, called the “Minimum Information Specification For In Situ Hybridization and Immunohistochemistry Experiments (MISFISHIE)”. It is modelled after the MIAME (Minimum Information About a Microarray Experiment) specification for microarray experiments. Data specifications like MIAME and MISFISHIE specify the information *content* without dictating a *format* for encoding that information. The MISFISHIE specification describes six types of information that should be provided for each experiment: Experimental Design, Biomaterials and Treatments, Reporters, Staining, Imaging Data, and Image Characterizations. This specification has benefited the consortium within which it was initially developed and is expected to benefit the wider research community. We welcome feedback from the scientific community to help improve our proposal.

Background

High-throughput analyses of gene expression in biological samples (*e.g.*, transcript abundance using microarrays or protein abundance using proteomics) often do not provide information about the cell types or spatial domains within tissues that express the genes of interest, and may not reveal dynamic or transient gene expression. Consequently, such analyses are often followed by experiments to confirm the location and degree of gene expression by specific cell types within the tissue by probing with specific reporters for the genes of interest. In addition, the wealth of clinical information associated with tissue samples in large collections all over the world provide a powerful tool to validate or expand the conclusions made using such high-throughput analysis of fresh samples.

However, it is often the case that studies that make use of *in situ* hybridization (ISH) and immunohistochemistry (IHC) staining, and/or their resulting images are presented without the information needed to interpret the images or the methodology that produced them. Furthermore, neither the reagents and methods used in the experiments, nor the results are easily searchable through current biomedical literature databases like PubMed. Since the interpretation of ISH and IHC stains could differ between observers, between different image analysis platforms and programs, and even between different sessions using the same image analysis platform and

program¹, communicating the methods and criteria used are critically important for teaching others and to permit critical evaluation of a published work.

Data annotation specifications that have been developed by the wider microarray community²⁻⁴ have begun to show benefits for the biomedical research community. First and foremost, the debate initiated by the proposal for specifications engaged many researchers, and the current specifications included the contributions of many different interests within the microarray data generating community. Common exchange formats and the willingness of researchers to put their data into the public domain upon publication have significantly increased the accessibility of data to all researchers. The open-source software and ontologies developed in conjunction with the data specifications resulted from the efforts of many different groups in the community. General discussion forums facilitated interaction between manufacturers and experimenters working towards development of the specifications for better experiments and better publications. Similar specifications are currently under development for other high-throughput technologies⁵⁻¹⁰.

Others have proposed data formats to better enable exchange of microscopy image data. For example, an XML data format specifically for tissue microarrays has been proposed¹¹. However, no minimum amount of information is specified, and users are free to include only as much information as they wish. Also available is Open Microscopy Environment (OME), which provides a flexible XML data format for storing and transmitting metadata for microscopy image datasets (<http://www.openmicroscopy.org/>). However, there is no comprehensive specification for facilitating the exchange of data from visual interpretation-based tissue protein and transcript abundance/localization experiments (hereafter referred to as ‘gene expression localization experiments’), such as *in situ* hybridization and immunohistochemistry.

Results and Discussion

To maximize the benefit of new gene expression localization experiments to the biomedical research community, we propose a minimum information specification, the “Minimum Information Specification For In Situ Hybridization and Immunohistochemistry Experiments (MISFISHIE)”. This specification provides guidelines for the minimum information that should be provided when publishing, making public, or exchanging results from visual interpretation-based tissue gene expression localization experiments such as ISH, IHC, lectin affinity histochemistry, and reporter gene constructs (e.g., green fluorescent protein [GFP], β -galactosidase). Compliance with this specification is expected to provide researchers at different laboratories with enough information to fully evaluate the data and to reproduce the experiment. Although MISFISHIE facilitates the identification of specific sources of variability, it cannot, and does not aim to, reduce this variability. However, if complete information, including raw image data, is always provided, the original interpretations may be re-evaluated by other researchers.

Modelled after the widely accepted MIAME specification for microarray experiments², MISFISHIE only prescribes the kind of information that should be provided. It does not include every parameter that could be specified about an

experiment, but rather broad categories of detail that should be addressed, relying on the data producers and reviewers to ensure that each section contains enough information for readers to be able to fully assess the validity of, and accurately reproduce the experiment described. Just as MIAME has been a guide to help authors provide enough information about a microarray experiment such that its interpretation could be verified or refuted¹², we hope that MISFISHIE will be used in the same way for gene expression localization experiments.

This specification does not dictate a specific format for reporting the information. We expect to develop a data model based on the concepts of MAGE-OM (MicroArray Gene Expression Object Model) and software based on the MAGEstk (MicroArray Gene Expression software tool kit) in the near future³. It is this model and the associated XML-based mark-up language that will provide a data format for archiving or transferring data. Since a major revision of MAGE-OM, the FuGE-OM (Functional Genomics Experiment Object Model)¹³, is currently being developed to accommodate data from other functional genomics experiments, it is likely that the MISFISHIE-derived object model will be an extension of FuGE-OM and not a separate construct. It is intended that MISFISHIE should function together with other technology-related specifications such as MIAME and MIAPE (Minimum Information About a Proteomics Experiment)¹⁴ to support functional genomics investigations. We anticipate that MISFISHIE will be integrated with other MGED (Microarray Gene Expression Data) Society standards¹⁵ through the Reporting Structure for Biological Investigation (RSBI) working group¹⁶. This is especially important since the goal of integrating different data types will most easily be realized when a common reporting structure is used. Separation of the minimum information specification and the data format is important because there should be scope for the provision of unlimited additional information beyond the minimum specification and encoding of incomplete information for optimal flexibility. Furthermore, broad acceptance of the minimum information required would greatly aid the design of a data model.

It has long been appreciated that improved standards for IHC are needed. However, standardization discussions have largely been focused on development of standardized technical protocols that might lead to more uniform staining¹⁷, or efforts towards reducing the subjectivity in interpretation of histological sections¹⁸. Here we do not attempt to endorse standardized methodologies or data interpretation, but rather seek to promote complete disclosure of the methodologies used so that experiments may be replicated by others employing the same procedures as the original investigators.

A set of guidelines specifically for tumor marker prognostic studies called REMARK¹⁹ has recently been established. REMARK encompasses the domain of outcome studies based on tumor markers of any kind, not just those of IHC. MISFISHIE encompasses the domain of studies employing IHC or ISH techniques; it may be a tumor marker study or a zebrafish embryo study. We believe that MISFISHIE presents a subset of guidelines applicable to nearly any IHC or ISH study regardless of context. We fully expect that specialized subdomains (such as clinical prognostic studies) will want to add applicable requirements for that subdomain.

While no accepted minimum specification for this type of data yet exists, there have been several efforts at organizing gene expression localization data in databases. Such

database designs provide a useful framework from which to build a specification. Two databases for the mouse research community, the Mouse Gene Expression Database (GXD)²⁰ and the Edinburgh Mouse Atlas Gene Expression (EMAGE) database²¹, have influenced the design of MISFISHIE. Mouse-specific fields in these databases were removed in favor of more organism-neutral ones. Several fields in these databases were deemed useful but are not part of a minimum requirement and, consequently, were not included. Also, in these databases many experiments that were entered by curators using the information provided in journal articles have empty fields because they had not been described in sufficient detail in the papers. Achieving MISFISHIE compliance in a publication will result in more complete reporting of experiments and, therefore, more reproducible experiments in these and other databases in the future. Although MISFISHIE is primarily designed as a specification for peer-reviewed journal articles, it will guide database development as well. The inclusion of specific experimental details, such as tissue type, reagents and methods, will allow investigators to find precedent for experiments they are considering more efficiently. For example, an investigator might be able to rapidly search all publications that reported immunoperoxidase localization of CD10 in the human prostate using the database and retrieve information on how the gene localization experiments were conducted.

This specification describes the type of information that should be provided for publication of gene expression localization experiments in six sections (Figure 1):

1. Experimental Design
2. Biomaterials (specimens) and Treatments (section or whole-mount preparation)
3. Reporters (probes or antibodies)
4. Staining
5. Imaging Data
6. Image Characterizations

The following description provides guidelines for ensuring that data are compliant with the specification. It is intended to be useful to researchers preparing to publish data as well as to manuscript reviewers and editors checking for MISFISHIE compliance. The use of ontologies, such as the MGED Ontology (MO)⁴ or Ontology for Biomedical Investigations (OBI; formerly named FuGO)²², facilitates computational searches of data and are therefore extremely advantageous as a source of descriptors. For terms outside the scope of OBI, such as those in anatomy, another appropriate ontology may be used. A good list of ontologies is maintained at the Ontologies for Biology Organization (OBO) web site <http://obo.sourceforge.net/>. Use of OBI and other ontologies will be especially important as MISFISHIE-supporting applications and databases are developed. Many of the terms used in this specification are already defined in OBI.

Experimental Design:

This section should contain information about the gene expression localization experiment as a whole including a brief description of the project, experimental factors, and the methods. For example, this would include the variables between the assays in the experiment, and how and where to get more information about the

experiment (web sites and contact persons). We propose that the following types of information be used to describe the overall design of an experiment:

- Experiment description: a short summary of the aims of the experiment.
- Assay type(s): e.g., immunohistochemistry, *in situ* hybridization, lectin affinity histochemistry, cell-lineage- or tissue-specific reporter expression.
- Experiment design type: e.g., is it a comparison of normal vs. diseased tissue, of multiple tissue/embryo specimens of similar type, of multiple probes/antibodies applied to the same tissue, or a localization screen, etc.? The MGED Ontology ExperimentDesignType has many entries categorizing design type.
- Experimental factors: the parameters or conditions that are tested, such as probe/antibody, disease state, genetic variation, structural unit, age, etc. Again, the MGED Ontology is a rich source of terms that can be used to describe the factors being tested.
- Total number of assays performed in the experiment: an assay is defined as one instance of a hybridization/stain of a single specimen with a single reporter. Thus, the result of a tissue microarray consisting of a 10 x 10 array of tissues would be counted as 100 assays. If replicates or reruns are a component of the experimental design, provide details that should include number of replicates per tissue, per reporter, etc.
- URL of any websites or database accession numbers (if available) pertinent to the experiment.
- Contact information for communicating with the experimenters.

Biomaterials (specimens) and Treatments (section or whole-mount preparation):

Describing specimens comprehensively is challenging, since they may have dozens or even hundreds of characteristics, especially for patient material when clinical information is available. The guiding principle in sample description is to supply enough information for an independent researcher to carry out a similar experiment. Characteristics that are known to differ among specimens should be provided with each specimen; while common attributes of all the specimens may be provided only once. The MISFISHE proposal lists characteristics of a biological sample that should be described:

- Origin of the biological specimens. Information required includes:
 - Attributes of the individual(s). The organism species must be named, preferably using the NCBI taxonomy, and for non-human organisms the strain and mutant alleles should be named according to the accepted standards for that organism. Additional attributes may include, but are not limited to, sex, age, developmental stage, genotype, phenotype.
 - Physiologic state of the individual(s) (normal vs. diseased).

- Relevant exogenous factors (e.g., treatment, special diet).
- Anatomic source of the tissue or cell sample.
- Provider of the specimens.

All information critical for other researchers to reproduce the biomaterials as closely as possible should be provided. The information is not limited to the above examples. Referencing an established ontology or controlled vocabulary for the terms used is *highly* encouraged. Ontologies and controlled vocabularies are available from many sources in a variety of formats, including on-line references and reference textbooks. Since we are still at an early stage in the development and widespread use of databases to store sample information, a standardized set of terms and a single, widely accepted ontology is not yet available. The rationale for providing specific structural detail is that the location of an object, such as a cell type that is being studied may correlate with expression of a specific gene by that cell type. Structural detail may be important not only for cases where gene expression is dependent on tissue handling (e.g., there is stronger labeling at the specimen edges), but also in cases where, even within a single microanatomical unit there is heterogeneity (e.g., in lung tumors, cell cycle regulatory genes are highest at the periphery)²³.

- Manner of preparation of the specimens for the study. Information required includes:
 - Nature of the specimens (e.g., whole tissue, whole mounts of tissue, tissue sections, thickness of sections, whole cells, or sections of cells).
 - Manner in which the specimens were prepared for the experiments (e.g., fixation with type of fixative and duration of fixation vs. fresh, non-fixed, non-frozen specimens or frozen specimens, sections mounted on slides versus sections floating in reagents).
 - Protocols used. Referencing previously published protocols is permissible if the protocols are appropriately detailed and were strictly followed.

Sensitivity of the immunoreaction of some gene products to fixation is exemplified by the observation that p27 was least frequent and least intense in prostate cancer cells that were farthest from the cut surface of a fixed tissue. These were the cells that are least rapidly fixed²⁴.

Reporters (probes or antibodies):

It is critical to provide full information about the reporters (probes, lectins, or antibodies) used, since these can differ in reactivity from lot to lot and manufacturer to manufacturer. A manufacturer's literature usually provides most of the needed information; however, the manufacturer's literature may not be permanent. For privately produced reporters, enough information needs to be provided so that another lab could produce the same compounds. MISFISHIE specifies several requirements

necessary to best describe the molecules used to label a tissue sample. It was noted in the review of this manuscript that thorough validation of reporters is very often poorly done in current literature. This specification does not at present require that researchers validate each reporter used in a particular way, but such validation is encouraged and should be reported when performed.

- Unambiguous genomic identification of each reporter:
 - For in situ hybridizations, provide the corresponding GenBank/EMBL/DDBJ accession number and, if applicable, the start and end nucleotide positions of the probe within that sequence. Also, provide the accession number version or database release version.
 - For antibodies, provide the protein identifier, including specific version information for the accession number or database release.
- Full sequence of each probe, or clone number of each antibody. For fluorescent protein experiments, the promoter sequence should be specified. In each case, provide the method by which the reporter was characterized.
 - If the sequence or clone number is not known, then the template or clone must be made publicly available. Provide specific details on how the template or clone may be obtained.
 - Some tissue localization experiments are based on the principle that the gene being localized is detected when the gene promoter activates a fluorescent protein reporter, such as GFP. In such experiments, the sequence of the reporter, i.e., GFP, is not important. Rather, the sequence of the promoter is critical and confers cell and tissue specificity to the reporter since the promoter is specific to that cell.
- Protocol(s) for how the reporters were designed and produced or the source from which they were obtained.
 - For reporters purchased from a company, the company name, address, catalogue number, and lot number should be provided.
 - For a custom-made antibody, the putative antigen and references to studies that characterize the sensitivity and specificity of the antibody in tissue immunostains.
- Additional attributes of the reporter:
 - For antibodies, the type of primary antibody (monoclonal or polyclonal), the immunoglobulin isotype, and the organism in which the antibody was generated.
 - For lectins, the full name (e.g., *Dolichos biflorus*), the source of the lectin (e.g., which company produced it), how it was detected (e.g., whether it was fluorescently labelled or biotinylated, with follow-up histochemical analysis), and how it was labelled (e.g., if the investigators labelled the lectin themselves the source of the reagents, the method and/or the labeling kit should be provided).

Staining:

The protocols used for staining vary considerably among experimenters. The merits of standardizing these protocols have been discussed extensively in the literature. This specification merely requires that the protocol used is provided and is sufficiently detailed that another researcher may follow it. The following types of information should be provided to adequately describe the staining protocols and parameters:

- Number of detectable reporters in the hybridization or stain (e.g., more than one for multiple-dye fluorescence microscopy) plus specific details about the detection method:
 - Detection reagent used (e.g., fluorophores used, enzyme-substrates, gold particles).
 - Source of the detection system plus sufficient detail to reproduce the reaction.
- Protocol used to produce the hybridization or immunostain. This should include a description of how the tissue (organism, organ, or section) was mounted onto the slide/substrate and treatments of the section, e.g., immunohistochemistry protocol inclusive of parameters such as buffer, temperature, post-wash conditions, etc. Referencing previously published protocols is permissible if the protocols are appropriately detailed and were strictly followed. Also include:
 - What steps, if any, were taken to decrease non-specific reaction product. For example, in immunoperoxidase experiments there might be pre-incubation of the specimen preparation with (a) albumin solution to block non-specific binding, (b) peroxide solution to block signal due to endogenous peroxidase.
 - Use of an antigen or gene product retrieval method.
- Information about assay controls: the nature of both positive and negative tissue and reporter controls (or state if controls were not performed). The same level of detail of the tissue controls should be reported as for the cells or tissues that are being studied. Optionally provide specificity reporter controls, such as competitive inhibition with either purified protein or peptide in immunohistochemistry.

Imaging Data:

Although the MIAME specification stops short of requiring microarray image data, we propose that MISFISHIE require that representative IHC or ISH images be provided since the interpretation of these images varies with the experience and training of the observer. While the images are not needed to facilitate reproducibility of an experiment, they greatly aid in the interpretation and analysis and in determining reasons for discordant results. Both positive and negative results should be reported;

this information is potentially useful for other work outside the scope of the reported experiment.

For several model organisms, there are already repositories for gene expression localization experiment images, including GXD²⁰ and EMAGE²¹ for mouse, ZFIN²⁵ for zebrafish, and others. However, for many organisms including human, there may not be such a dedicated database. It would be of tremendous value to the research community to have a general, organism-independent database for archiving gene expression localization experiment images. Such an archive could provide examples of tissue localization studies, and could be a reference site for investigators who want to verify the tissue localizations of reporter reagents they are considering using. More importantly, a general-purpose repository to which researchers could submit their images for permanent storage with accession numbers for publications would be very valuable for facilitating MISFISHIE compliance and in realizing the full value of these data for future research. MorphBank (<http://www.morphbank.net>) is an available general purpose image repository for biological research. BioImage is an image repository under construction at <http://www.bioimage.org/>²⁶.

The MISFISHIE specification suggests that the following information should be provided:

- Digital images for each assay included in the study should be digitally available for download without additional charge. The images should be of sufficient resolution to allow independent characterization, and provided in a standard file format (e.g., JPEG, PNG, GIF, TIFF). The images should be named or tagged with the reporter and specimen that they represent.
- Detection method by which hybridization or staining is observed (e.g., for each channel a fluorescent wavelength if multiple reporters are used). If the detection method is the same for all images, it need only be mentioned once.
- Images for the controls are not required, although may optionally be provided.

Image Characterizations:

The results as interpreted by the original researchers should be reported in a clearly articulated, concise and consistent manner. This permits reviewers to ensure that the characterizations are consistent with and representative of the data, and that the conclusions are reasonable. The characterizations should also be provided in such a way that they can be easily stored in a database, queried, and compared with other expression data.

The types of characterization recorded can vary depending on the experimental design. The following guidelines specify a minimum set of characterization features. Additional characterization of the images as required by the experimental design could also be provided.

- Ontology entries, including reference to the ontology (e.g., refs. ²⁷⁻³⁰, note that some ontologies, such as SNOMED CT and NIH/NLM's Unified Medical Language System (UMLS), may contain licensing restrictions that makes them unavailable to some or limits the use of the terms; a MISFISHIE-compliant document that contains SNOMED CT entries or some UMLS entries may not be legally redistributable³¹), terms, accession numbers, or terms and definitions if sufficient detail cannot be found in an existing ontology for individual structural units used for classification. Structural units could be an organ, tissue, cell, subcellular component, etc. Note that only the structural units relevant to the experiment need to be listed and characterized. It is not necessary to list (and characterize) structural units visible in the assays or slides but not relevant to the experiment or report.
- Intensity scale, ideally choosing one from the MGED Ontology. For example, a three-level scale of present, absent, or equivocal might be appropriate for evaluating IHC stains. However, any scale that the investigators feel is appropriate may be used as long as each gradation of intensity in the scale is defined in a manner that enables an independent investigator to understand or apply the same criteria.
- Per each structural unit (relevant to the experiment) in each assay (or in each image), provide:
 - Staining intensity, or the fraction of the structural unit's population exhibiting each intensity (see example below).
 - Other optional annotations/characterizations of the structural unit, e.g., feature density, qualitative characteristics or spatial distribution of the structural unit or staining. The use of referenced ontology terms is encouraged.

Both positive and negative calls of staining relevant to the experiment should be reported. It is quite useful to provide negative expression results; it is understood that a negative result is actually an upper limit to the expression level, where the limit is usually not well known. If some structural units cannot be characterized for some reporters, corresponding calls may be null.

For example:

Luminal epithelial cell: present
 Basal epithelial cell: absent
 etc.

is sufficient; or, when appropriate for the type of analysis being done, more detail:

Luminal epithelial cell: 90% present, 10% equivocal, 0% absent
 Basal epithelial cell: 10% present, 10% equivocal, 80% absent
 etc.

Unless only a few expression calls are presented, it is clearest if the calls are presented in tabular form, either within the manuscript or as supplemental material, as appropriate.

- Optionally as a best practice, the protocol for the characterization and information about the basic technique for characterizing the assays. For

example, this information may include how many observers performed the characterizations, whether the characterizations were performed from the images themselves or visually through the instrument, any exceptions or assumptions made in characterizing the data, etc. We refer to one example of a well-described characterization protocol³². We also note that it has been reported that performing the characterization from digital images has advantages in terms of replication, decreased intraobserver and interobserver variability³³.

Some examples of real experimental data annotated according to MISFISHIE are posted at the MISFISHIE web site, available as a link from the MGED workgroup web page <http://www.mged.org/Workgroups/>. We also provide an abbreviated checklist (Figure 2) to aid in assessing MISFISHIE compliance. It should be use in conjunction with the full description, not in place of it. A printable version is supplied as Additional File 1.

Survey of the recent literature

To assess how the MISFISHIE specification compares with what appears to be standard practice for publication today, a selection of articles reporting on IHC or ISH from the last five years were assessed for compliance with the six sections of the MISFISHIE specification. Three articles³⁴⁻³⁶ were assessed and discussed by all ten *ad hoc* reviewers so that inter-reviewer variability could be minimized. Another 29 articles³⁷⁻⁶⁴ were assigned to individual reviewers for assessment. Each reviewer assessed each of his or her assigned articles in the context of a scenario of a journal referee reviewing a submitted article. As part of the review, the MISFISHIE compliance checklist (see Additional File 1) was completed by the reviewer as if it were the journal's policy to require MISFISHIE compliance.

Compliance for each MISFISHIE subsection was rated by the reviewers on the scale of 0 to 10, where a 10 indicates that the authors provided all information that the reviewer needed to understand or reproduce the experiment without needing to make any assumptions. Scores lower than 10 correspond to how incomplete the information was that the reviewer thought necessary to understand or reproduce the work. Scores of 8 and 9 were considered a low pass; the reviewer could reproduce the experiment although with a few assumptions. It was therefore possible for a paper to leave out a few details that the reviewers deemed ought to have been provided, but still pass. Compliance with each section was somewhat subjective as the strictness of each reviewer was not uniform, as would presumably be the case for *bona fide* journal reviewers. Therefore, the MISFISHIE specification itself is subject to individual interpretation. Since this cannot be avoided, we hope that the checklist will minimize subjectivity.

This exercise not only proved useful in testing the proposed MISFISHIE specification, but also allowed us to determine if any section seemed too onerous a requirement. Of the 32 papers assessed, only four (13%) were deemed MISFISHIE compliant in all six sections. An additional 28% were out of compliance with only one section, and 31% did not comply in two sections. The review considered that more than 90% of the papers were compliant with MISFISHIE sections 1 and 2

(Experimental Design; Biomaterials and Treatments). Compliance for sections 3 and 4 (Reporters and Staining) was about 75%. Section 5 (Imaging Data) proved to be the most troublesome, with only 16% of the articles compliant. Finally, about 47% complied with section 6 (Image Characterizations). These results are summarized in Table 1.

Although few of the surveyed articles complied fully, the reviewers felt that the majority of non-compliant papers would require only modest additions to become compliant, with the possible exception of section 5. This section requires that at least one representative image of each assay be electronically available. This may be within a model organism database, a generic image database, a journal's supplemental data web site, or even the author's web site, although the latter is the least preferable. It is not necessary for all images to be reproduced within the manuscript itself. One might feel that making all images accessible to others can be unduly burdensome. However, we feel that since image interpretation is variable, it is necessary that the original images be made available in a digital format for subsequent review, ideally in a centrally-managed public repository. Some model organism databases already provide such a facility. MorphBank provides an example of a general-purpose image repository for any organism, although it does not appear to be well suited to store the accompanying characterizations in an easily queryable format.

We provide as one example of a paper that was deemed MISFISHIE compliant the work of Santagata et al.⁵⁹ Our review of this article concluded that it provides sufficient detail for all MISFISHIE sections; all images used for the study are available at their own website.

Conclusions

This specification was jointly developed by members of the NIH/NIDDK Stem Cell Genome Anatomy Projects consortium to facilitate data sharing within the consortium. After use and refinement within the consortium, and based on discussions with additional members of the larger research community, we offer this specification, published here as MISFISHIE version 1.0 as a proposal to the whole research community. The history of the creation of MISFISHIE and the lessons learned from it⁶⁵ may be helpful for others aiming to create a similar specification for other data types.

We expect that MISFISHIE will undergo updates, leading to future editions, as other localization methods, such as DNA *in situ* hybridization experiments to chromosomes, are implemented and the need for a specification is expressed. The eventual accepted specification cannot be dictated, but rather must be achieved through discussion and consensus. Suggestions from the community are actively encouraged and will be collected and folded into an eventual second release, published at the MISFISHIE area of the MGED website: <http://www.mged.org/Workgroups/MISFISHIE/>. Comments may be addressed to the email distribution list dedicated to discussion about MISFISHIE: mged-misfishie@lists.sourceforge.net. We note that there is still considerable room for researching the scientific best practices for performing and reporting these types of studies. We have attempted here to define a minimum set of information and have

provided a few optional best practices that were deemed not quite appropriate as a requirement for all publications.

After a suitable period of dialog and revision by the community, and should the community accept the final proposal, we would encourage reviewers, journal editors and funding agencies to promote compliance with MISFISHIE for all studies that report gene expression localization data so that all published data and resulting conclusions may be correctly interpreted, and that independent investigators would have the necessary information that would enable them to repeat the experiment.

Our survey of recent articles indicated that only about 15% of published works are fully compliant with this specification, and most fail by not making images of assays used in the study digitally accessible to the research community. Most of the surveyed papers could be brought into compliance by uploading the images into a repository and adding fewer than a dozen additional sentences of description. If article length constraint hinders full MISFISHIE compliance, it would be encouraged that the information be provided in supplemental material.

Several of the model organism databases are already able to accept and archive the results from a publication that provides all information that MISFISHIE specifies. We highly encourage authors to submit their data to these databases via the provided database submission process upon submission of the article.

List of Abbreviations

EMAGE: Edinburgh Mouse Atlas Gene Expression database
(<http://genex.hgu.mrc.ac.uk/>)

FuGE-OM: Functional Genomics Experiment Object Model

FuGO: Functional Genomics Ontology (renamed OBI in Oct 2006)

GFP: green fluorescent protein

GXD: Gene Expression Database (<http://www.informatics.jax.org/>)

IHC: immunohistochemistry

ISH: *in situ* hybridization

MAGE-OM/ML: MicroArray Gene Expression Object Model/ Markup Language

MGED: Microarray Gene Expression Data Society (<http://www.mged.org/>)

MIAME: Minimum Information About a Microarray Experiment

MIAPE: Minimum Information About a Proteomics Experiment

MISFISHIE: Minimum Information Specification For In Situ Hybridization and Immunohistochemistry Experiments

MO: MGED Ontology

NIDDK: National Institute of Diabetes & Digestive & Kidney Diseases

NIH: National Institutes of Health

NLM: National Library of Medicine

OBI: Ontology for Biomedical Investigations (formerly FuGO)

OME: Open Microscopy Environment (<http://www.openmicroscopy.org/>)

PEDRo: Proteomics Experiment Data Repository (<http://pedro.man.ac.uk/>)

RSBI: Reporting Structure for Biological Investigation

UMLS: Unified Medical Language System

XML: Extensible Markup Language

Acknowledgements

We thank Rachel Drysdale, Lillian Eichner, Mervi Heiskanen, and Monte Westerfield for comments and discussions during the preparation of the MISFISHIE specification, and Christine Emswiler for assistance with the figures. This work was funded in part with support from the National Institute of Diabetes & Digestive & Kidney Diseases, National Institutes of Health, to members of the Stem Cell Genome Anatomy Project Consortium, including. DK63483 to Jeff Gordon (Washington University in St. Louis), DK63481 to Ihor Lemischka (Princeton University), DK63400 to Melissa Little (University of Queensland), DK63630 to Alvin Liu (University of Washington), and DK63328 to Len Zon (Children's Hospital Boston).

References

1. True, L.D. Quantitative immunohistochemistry: a new tool for surgical pathology? *Am J Clin Pathol* **90**, 324-325 (1988).
2. Brazma, A. et al. Minimum information about a microarray experiment (MIAME)-toward standards for microarray data. *Nat Genet* **29**, 365-371 (2001).
3. Spellman, P.T. et al. Design and Implementation of Microarray Gene Expression Markup Language (MAGE-ML). *Genome Biology* **3**, RESEARCH0046 (2002).
4. Stoeckert, C.J. & Parkinson, H. The MGED ontology: A Framework for Describing Functional Genomics Experiments - <http://mged.sourceforge.net/ontologies/MGEDontology.php>. *Comparative and Functional Genomics* **4**, 127-132 (2003).
5. Taylor, C.F. et al. A Systematic Approach to Modeling Capturing and Disseminating Proteomics Experimental Data. *Nature Biotechnology* **21**, 247 (2003).
6. Garwood, K. et al. PEDRo: A database for storing, searching and disseminating experimental proteomics data. *BMC Genomics* **5**, 68 (2004).
7. Jones, A., Hunt, E., Wastling, J.M., Pizarro, A. & Stoeckert, C.J., Jr. An object model and database for functional genomics. *Bioinformatics* **20**, 1583-1590 (2004).
8. Xirasagar, S. et al. CEBS object model for systems biology data, SysBio-OM. *Bioinformatics* **20**, 2004-2015 (2004).
9. Jenkins, H. et al. A proposed framework for the description of plant metabolomics experiments and their results. *Nat Biotechnol* **22**, 1601-1606 (2004).
10. Lindon, J.C. et al. Summary recommendations for standardization and reporting of metabolic analyses. *Nature Biotechnology* **23**, 833-838 (2005).
11. Berman, J.J., Edgerton, M.E. & Friedman, B.A. The tissue microarray data exchange specification: a community-based, open source tool for sharing tissue microarray data. *BMC Med Inform Decis Mak* **3**, 5 (2003).
12. Stoeckert, C.J., Quackenbush, J., Brazma, A. & Ball, C.A. Minimum information about a functional genomics experiment: the state of microarray

- standards and their extension to other technologies. *Drug Discovery Today: TARGETS* **3**, 159-164 (2004).
13. Jones, A.R., Pizarro, A., Spellman, P., Miller, M. & group, T.F.w. FuGE: Functional Genomics Experiment Object Model. *OMICS: A Journal of Integrative Biology* **10**, 179-184 (2006).
 14. Taylor, C.F. et al. HUPO - Proteomics Standards Initiative (PSI). *IOMICS: A Journal of Integrative Biology* **in press** (2006).
 15. Ball, C.A. & Brazma, A. MGED Standards. *OMICS: A Journal of Integrative Biology* **10**, 138-144 (2006).
 16. Sansone, S.-A. et al. A Strategy Capitalizing on Synergies: The Reporting Structure for Biological Investigation (RSBI) Working Group. *OMICS: A Journal of Integrative Biology* **10**, 164-171 (2006).
 17. Swanson, P.E. Methodologic Standardization in Immunohistochemistry: A Doorway Opens. *Applied Immunohistochemistry* **1**, 229-231 (1993).
 18. Taylor, C.R. An exaltation of experts: concerted efforts in the standardization of immunohistochemistry. *Hum Pathol* **25**, 2-11 (1994).
 19. McShane, L.M. et al. Reporting recommendations for tumor marker prognostic studies (REMARK). *J Natl Cancer Inst* **97**, 1180-1184 (2005).
 20. Hill, D.P. et al. The mouse Gene Expression Database (GXD): updates and enhancements. *Nucleic Acids Res* **32 Database issue**, D568-571 (2004).
 21. Baldock, R.A. et al. EMAP and EMAGE: a framework for understanding spatially organized data. *Neuroinformatics* **1**, 309-325 (2003).
 22. Whetzel, P.L. et al. Development of FuGO – an Ontology for Functional Genomics Experiments. *OMICS: A Journal of Integrative Biology* **10**, 199-204 (2006).
 23. Dobashi, Y. et al. Active cyclin A-CDK2 complex, a possible critical factor for cell proliferation in human primary lung carcinomas. *Am J Pathol* **153**, 963-972 (1998).
 24. De Marzo, A.M., Fedor, H.H., Gage, W.R. & Rubin, M.A. Inadequate formalin fixation decreases reliability of p27 immunohistochemical staining: probing optimal fixation time using high-density tissue microarrays. *Hum Pathol* **33**, 756-760 (2002).
 25. Sprague, J. et al. The Zebrafish Information Network (ZFIN): the zebrafish model organism database. *Nucleic Acids Res* **31**, 241-243 (2003).
 26. Carazo, J.M. & Stelzer, E.H. The BioImage Database Project: organizing multidimensional biological images in an object-relational database. *J Struct Biol* **125**, 97-102 (1999).
 27. Rosse, C. & Mejino, J.L., Jr. A reference ontology for biomedical informatics: the Foundational Model of Anatomy. *J Biomed Inform* **36**, 478-500 (2003).
 28. Bard, J., Rhee, S.Y. & Ashburner, M. An ontology for cell types. *Genome Biol* **6**, R21 (2005).
 29. Bard, J.L. et al. An internet-accessible database of mouse developmental anatomy based on a systematic nomenclature. *Mech Dev* **74**, 111-120 (1998).
 30. Hayamizu, T.F., Mangan, M., Corradi, J.P., Kadin, J.A. & Ringwald, M. The Adult Mouse Anatomical Dictionary: a tool for annotating and integrating data. *Genome Biol* **6**, R29 (2005).
 31. Berman, J.J. A tool for sharing annotated research data: the "Category 0" UMLS (Unified Medical Language System) vocabularies. *BMC Med Inform Decis Mak* **3**, 6 (2003).

32. Abd El-Rehim, D.M. et al. Expression of luminal and basal cytokeratins in human breast carcinoma. *J Pathol* **203**, 661-671 (2004).
33. Bova, G.S. et al. Web-based tissue microarray image data analysis: initial validation testing through prostate cancer Gleason grading. *Hum Pathol* **32**, 417-427 (2001).
34. Liu, A.Y. & True, L.D. Characterization of prostate cell types by CD cell surface molecules. *Am J Pathol* **160**, 37-43 (2002).
35. Kernek, K.M. et al. Fluorescence in situ hybridization analysis of chromosome 12p in paraffin-embedded tissue is useful for establishing germ cell origin of metastatic tumors. *Mod Pathol* **17**, 1309-1313 (2004).
36. McKenney, J.K. et al. Basal cell proliferations of the prostate other than usual basal cell hyperplasia: a clinicopathologic study of 23 cases, including four carcinomas, with a proposed classification. *Am J Surg Pathol* **28**, 1289-1298 (2004).
37. Amara, N. et al. Prostate stem cell antigen is overexpressed in human transitional cell carcinoma. *Cancer Res* **61**, 4660-4665 (2001).
38. Ayala, G. et al. High levels of phosphorylated form of Akt-1 in prostate cancer and non-neoplastic prostate tissues are strong predictors of biochemical recurrence. *Clin Cancer Res* **10**, 6572-6578 (2004).
39. Bart, J. et al. The distribution of drug-efflux pumps, P-gp, BCRP, MRP1 and MRP2, in the normal blood-testis barrier and in primary testicular tumours. *Eur J Cancer* **40**, 2064-2070 (2004).
40. Browne, T.J. et al. Prospective evaluation of AMACR (P504S) and basal cell markers in the assessment of routine prostate needle biopsy specimens. *Hum Pathol* **35**, 1462-1468 (2004).
41. Chen, D. et al. Syndecan-1 expression in locally invasive and metastatic prostate cancer. *Urology* **63**, 402-407 (2004).
42. Clayton, H., Titley, I. & Vivanco, M. Growth and differentiation of progenitor/stem cells derived from the human mammary gland. *Exp Cell Res* **297**, 444-460 (2004).
43. Cooray, H.C., Blackmore, C.G., Maskell, L. & Barrand, M.A. Localisation of breast cancer resistance protein in microvessel endothelium of human brain. *Neuroreport* **13**, 2059-2063 (2002).
44. Giangreco, A., Shen, H., Reynolds, S.D. & Stripp, B.R. Molecular phenotype of airway side population cells. *Am J Physiol Lung Cell Mol Physiol* **286**, L624-630 (2004).
45. Gmyrek, G.A. et al. Normal and malignant prostate epithelial cells differ in their response to hepatocyte growth factor/scatter factor. *Am J Pathol* **159**, 579-590 (2001).
46. Hwang, J.H. et al. Isolation of muscle derived stem cells from rat and its smooth muscle differentiation [corrected]. *Mol Cells* **17**, 57-61 (2004).
47. Jonker, J.W. et al. The breast cancer resistance protein BCRP (ABCG2) concentrates drugs and carcinogenic xenotoxins into milk. *Nat Med* **11**, 127-129 (2005).
48. Knudsen, B.S. et al. High expression of the Met receptor in prostate cancer metastasis to bone. *Urology* **60**, 1113-1117 (2002).
49. Larkin, A. et al. Investigation of MRP-1 protein and MDR-1 P-glycoprotein expression in invasive breast cancer: a prognostic study. *Int J Cancer* **112**, 286-294 (2004).

50. Lee, K., Klein-Szanto, A.J. & Kruh, G.D. Analysis of the MRP4 drug resistance profile in transfected NIH3T3 cells. *J Natl Cancer Inst* **92**, 1934-1940 (2000).
51. Li, R. et al. High level of androgen receptor is associated with aggressive clinicopathologic features and decreased biochemical recurrence-free survival in prostate: cancer patients treated with radical prostatectomy. *Am J Surg Pathol* **28**, 928-934 (2004).
52. Martin, C.M. et al. Persistent expression of the ATP-binding cassette transporter, Abcg2, identifies cardiac SP cells in the developing and adult heart. *Dev Biol* **265**, 262-275 (2004).
53. Martin, M.J., Muotri, A., Gage, F. & Varki, A. Human embryonic stem cells express an immunogenic nonhuman sialic acid. *Nat Med* **11**, 228-232 (2005).
54. Master, V.A., Wei, G., Liu, W. & Baskin, L.S. Urothelium facilitates the recruitment and trans-differentiation of fibroblasts into smooth muscle in acellular matrix. *J Urol* **170**, 1628-1632 (2003).
55. Piotrowska, A.P. et al. Alterations in smooth muscle contractile and cytoskeleton proteins and interstitial cells of Cajal in megacystis microcolon intestinal hypoperistalsis syndrome. *J Pediatr Surg* **38**, 749-755 (2003).
56. Ricciardelli, C. et al. Androgen receptor levels in prostate cancer epithelial and peritumoral stromal cells identify non-organ confined disease. *Prostate* **63**, 19-28 (2005).
57. Roudier, M.P. et al. Phenotypic heterogeneity of end-stage prostate carcinoma metastatic to bone. *Hum Pathol* **34**, 646-653 (2003).
58. Rubin, M.A. et al. Quantitative determination of expression of the prostate cancer protein alpha-methylacyl-CoA racemase using automated quantitative analysis (AQUA): a novel paradigm for automated and continuous biomarker measurements. *Am J Pathol* **164**, 831-840 (2004).
59. Santagata, S. et al. JAGGED1 expression is associated with prostate cancer metastasis and recurrence. *Cancer Res* **64**, 6854-6857 (2004).
60. Scotlandi, K. et al. C-kit receptor expression in Ewing's sarcoma: lack of prognostic value but therapeutic targeting opportunities in appropriate conditions. *J Clin Oncol* **21**, 1952-1960 (2003).
61. Shah, R.B. et al. Androgen-independent prostate cancer is a heterogeneous group of diseases: lessons from a rapid autopsy program. *Cancer Res* **64**, 9209-9216 (2004).
62. St Croix, B. et al. Genes expressed in human tumor endothelium. *Science* **289**, 1197-1202 (2000).
63. Wang, Z. et al. Expression of the human cachexia-associated protein (HCAP) in prostate cancer and in a prostate cancer animal model of cachexia. *Int J Cancer* **105**, 123-129 (2003).
64. Zhigang, Z. & Wenly, S. Prostate stem cell antigen (PSCA) expression in human prostate cancer tissues: implications for prostate carcinogenesis and progression of prostate cancer. *Jpn J Clin Oncol* **34**, 414-419 (2004).
65. Deutsch, E.W. et al. Development of the Minimum Information Specification For In Situ Hybridization and Immunohistochemistry Experiments (MISFISHIE). *OMICS: A Journal of Integrative Biology* **10**, 205-208 (2006).

Figures

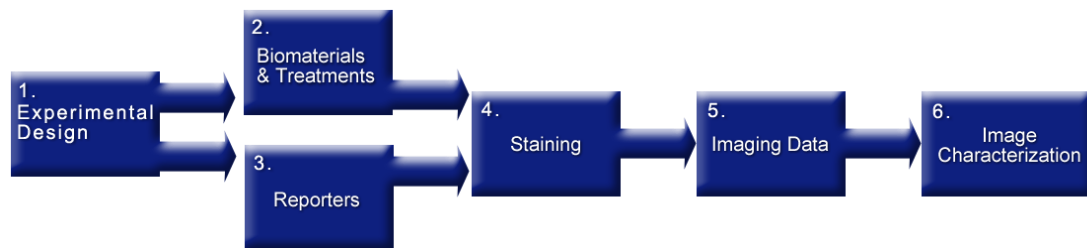


Figure 1: The six sections of the MISFISHIE specification.

Reviewer: _____ Date: _____

Paper Title: _____

Author List: _____

Jrnl, Yr, Pg: _____

The following checklist is an convenient abbreviated form of the full specification, and should be used in conjunction with the full specification to insure full compliance. The primary test for each line item is to judge whether enough information has been supplied to evaluate or reproduce the experiment without significant ambiguities.

Checklist:

	P/F	P/F	Subsection	Additional Comments
▶	Yellow	Black	1. ExperimentDesign	
			Experiment Description	
			Assay type(s) (IHC, ISH, GFP, etc.)	
			Experimental design (multiple reporter survey, specimen variation)	
			Experimental factors (variables in assays like reporter or specimen, etc.)	
			Total number of assays performed	
		Pink	(optional) URL for more information	
			Contact Information	
▶	Yellow	Black	2. BioMaterials & Treatments	
			Attributes of the individual (e.g., organism, sex, strain, line, dev stage, age, etc.)	
			Physiologic state (e.g., normal vs. disease)	
			Relevant exogenous factors (e.g., treatment, special diet, etc.)	
			Anatomic source of specimens	
			Provider of the specimens	
			Assay preparation Protocol (enough to reproduce?)	
▶	Yellow	Black	3. Reporter (Probe or Antibody) Information	
			Unambiguous reporter identification, ideally genomic	
			Full sequence or clone id of the reporters	
			Protocol for obtaining exact reporter (purchase from..., create, etc.)	
			Other important attributes (e.g., mono- or polyclonal, gen organism, etc.)	
▶	Yellow	Black	4. Staining Protocols & Parameters	
			Detection Method (number of reporters, det reagent & systems)	
			Staining protocol (enough to reproduce?)	
			Details about positive and negative controls	
▶	Yellow	Black	5. Imaging Data and parameters	
			The digital images for each assay (can download to your computer and explore?)	
		Pink	(optional) Imaging acquisition protocol	
▶	Yellow	Black	6. Image Characterizations	
			Definition of structural units (from ontology or manual definition)	
			Definition of intensity scale	
			Characterization of results in tabular form (digital or printed)	
		Pink	(optional) Characterization protocol	
▶	Yellow		Overall MISFISHIE Compliance (Any one section F is F)	

Comments:

Figure 2: An abbreviated checklist for the full MISFISHIE specification. This checklist should be used in conjunction with the full specification, not instead of it.

Tables

N	Percent	Statistic
32	100%	Number of articles assessed for compliance
4	13%	Number of articles considered to be fully MISFISHIE compliant
9	28%	Number of articles for which MISFISHIE information is missing for one section
10	31%	Number of articles for which MISFISHIE information is missing for two sections
6	19%	Number of articles for which MISFISHIE information is missing for more than two sections
31	97%	Number of articles that meet the data content requirements for section 1
29	91%	Number of articles that meet the data content requirements for section 2
24	75%	Number of articles that meet the data content requirements for section 3
24	75%	Number of articles that meet the data content requirements for section 4
5	16%	Number of articles that meet the data content requirements for section 5
15	47%	Number of articles that meet the data content requirements for section 6

Table 1: Summary of statistics from the MISFISHIE assessment survey of a cohort of selected current literature.

Additional files

Additional file 1: A simplified, printable, 1-page MISFISHIE checklist:
http://scgap.systemsbiology.net/standards/misfishie/MISFISHIE_abbrev_checklist.xls