# Redacting PHI in Neurological Images using XNAT

Alex Barclay,[1,2] Nakeisha Schimke,[2] John Hale [2]

[1] *Laureate Institute for Brain Research, Tulsa, Oklahoma*
[2] *Institute of Bioinformatics and Computational Biology, University of Tulsa, Tulsa, Oklahoma*

**Introduction**

In neuroscience, the potential for large scale collaboration and data sharing is undermined by concerns over the management of personal health information (PHI) in neuroimagery data sets. In particular, HIPAA and HITECH rules mandate substantial measures for the preservation of subject privacy, and for the researcher, the management of information privacy is a burdensome task. The potential for sharing and meta analysis of neuroimagery in research goes unfulfilled for two reasons: (1) sanitization is hampered by the lack of tools to systematically expunge PHI from DICOM images, and (2) an established workflow to integrate sanitization into the data sharing process remains conspicuously absent.

As the demand for medical imaging grows, it also increases in number of modalities, resolution of the images collected, and complexity of the data storage format. The amount of data generated from large-scale neuroimaging research studies and the necessary mandatory security and privacy filters make inter-organizational collaboration a daunting task. Current methods exist to redact, or verifiably sanitize, sensitive information for legal production and classified data; however, these methods are incomplete and only address small areas of a much larger problem.

In the US, a patient's right to privacy is mandated under the Health Insurance Portability and Accountability Act (HIPAA), which provides guidance for what patient data is considered protected health information (PHI). There are sixteen well-defined identifiers, including patient name, patient unique identifier, and geography. These can be thought of as keywords (ex. "Bob", "1234", and "1234 Main Street") and can be generally be found and isolated in legal redaction tools easily; however, no such software exists for the medical images themselves. There are two ambiguous identifiers specified within HIPAA: (1) full face photographic images and any comparable images, and (2) any other unique identifying number, characteristic, or code. These requirements necessitate a higher level of understanding of the presentation of the underlying data structure for correct and complete data removal, particularly when high resolution medical images may themselves be considered PHI.

To solve this problem, data must logically be deconstructed and considered as different layers in the hardware, operating system, software, and data analysis stack. This process provides a holistic framework for evaluating the current state of tools and techniques for each layer. This is different from even the best current legal redaction techniques that do not take into account the multiple levels, structure, or presentation of data. These restrictions allow the possibility of nefarious radiology to reconstruct facial features or other potentially identifiable neurocharacteristics. It is important that any redaction process is complete, and that all identifiers are removed from each layer. Otherwise, patient privacy cannot be guaranteed, and a potential breech of patient information can occur.

Through a comprehensive approach that considers data at every level in the stack, we create a comprehensive and practical infrastructure for managing PII in neuroimagery datasets, relieving the burden of the investigator from technical aspects of data sanitization and redaction. In doing so, we can remove substantial obstacles to large-scale collaboration and data sharing in neuroscience.

**Workflow**

When a patient is scanned for a study, the scanner records the data to a hard drive. This image is stored as a DICOM file, which contains the pixel data, patient identifying data, study parameters, and acquisition data. The image, before being shared, must be redacted. This can be considered a two-pass redaction, the first eliminating of the 16 well defined identifiers and the second the more ambiguous images. Since the first pass is fairly standard, it typically will not need to be verified, but we store the original unredacted form as well as the sanitized data, maintaining a mapping of original to redacted identifiers. This allows persistent pseudonymous identifiers that allow patients to be seamlessly tracked through multiple scans. The second pass then removes any facial features from this semi-anonymized version. This process is, like the HIPAA requirements it meets, more ambiguous, since it is not always agreed upon which, if any, image data should be removed. It, therefore, requires some user intervention to verify the data, and may require more than one method or parameter set. For this reason, we must maintain the original image. After both passes are complete, the image is ready for analysis and collaboration.

**Example Data Stack**
This begins at the simplest and lowest level in the data stack with the raw neuroimage data on the physical storage media (ex. hard drive) in its binary form as a stream of data. In the next layer, the file system, this is represented as a file pointing to an address space on disk along with associated metadata. The standardized file format for medical imaging is DICOM (Digital Imaging and Communications in Medicine). It contains both metadata and pixel data in a single file. The image pixel data can be encapsulated for compression purposes, which we represent as its own layer. The logical layer is the visual presentation of data.

*Storage Device*
A hard drive derives its value from the interpretation of data it stores; there are few places on a hard drive that have meaning on their own. When performing low level redaction on any of the data layers described above, the typical tool directs the hard drive to write low level data streams in the form of tokens or dummy data. There are mature methods to do this, and this does not have to be re-implemented.

*File System*
To comply with HIPAA, the file system must to be sanitized if it contains file structure layout including PHI information (folders, etc). This information can be redacted by a bitwise sliding window scan of PHI keywords and hash token replacement. Care must be taken to ensure file system integrity and preserve structure or keywords when transferring low level copies of data. Tools for this procedure are not satisfactory but can be easily implemented by building upon file system debugging tools such as debugfs for Linux.

*File Format & Encapsulated data*
Rather than a true header format, the DICOM standard approaches data as a list of attributes so that metadata, such as patient name and identifier, are part of the DICOM file. Pixel data is also considered an attribute, allowing the DICOM standard to encapsulate other formats and leverage existing compression techniques such as JPEG or run length encoding. This poses a challenge, since we must anticipate the use of numerous compression techniques and formats when considering a DICOM image. Since study metadata are also considered attributes, these are relatively straightforward to redact; however, care must be taken to not corrupt the final data, when these attributes are removed or replaced with pseudonymous tokens. This technique is also applied to file system metadata presented as file forks or extended attributes. There are several tools that currently perform this task sufficiently, and our tool will implement or extend a method similar to these existing tools to sanitize patient data.

*Logical Image*
The logical representation is a medical image, such as an MRI or CT scan. Complete images are typically three-dimensional, though we also consider each two-dimensional cross-section an image. While the image may not be immediately identifiable with respect to an individual, it may contain enough facial features to reconstruct the patient's face. This is considered PHI and should be redacted before it is shared. Single 2D slices, however, do not contain enough data to identify the patient and therefore do not need to be removed when stored individually. While there are tools currently in place to remove non-brain matter from a neurological image, our tool streamlines this process to make it as transparent as possible to the end user. This requires some user intervention to verify that the redacted image still contains a sufficient volume for analysis.

**Project Goal**
To satisfy HIPAA requirements while maintaining ease of use, we offer an automated approach to redacting patient data for use with the Extensible Neuroimage Archive Toolkit (XNAT) [7] as a repository, folding neatly within existing XNAT operational workflows and processes. This project has four primary advantages: (1) it provides a privacy map for linking redacted identifiers across institutional boundaries, (2) as part of the processing pipeline it is a natural part of a collaborative workflow, (3) it offers an unencumbered digital redaction building block, and (4) it supports verifiable redaction. The redaction effort consists of three logical components: the redaction engine, the Privacy Map, and the XNAT translator. The redaction engine contains both DICOM image redaction procedures and a low-level bitwise verification mechanism to fully remove embedded PHI. The privacy map is an encrypted database containing the link between the original images and their redacted counterparts, accessible by an API to ensure the integrity of the databases by exposing a non-harmful subset of commands to the redaction engine. It provides consistency between research subjects of the redacted ID to subject PHI, preventing statistical skew due to data duplication. The final component creates new XNAT identifiers and transmits redacted data.

By addressing redaction at multiple layers of the data stack, we can implement a comprehensive workflow that provides complete redaction to maintain the highest level of patient privacy.

**References**

[1] American College of Radiology, A. C. (2008), *Digital Imaging and Communications in Medicine (DICOM)*, National Electrical Manufacturers Association.

[2] Architectures and Applications Division of the Systems & Network Attack Center (2005),'Redacting with Confidence: How to Safely Publish Sanitized Reports Converted From Word to PDF', Technical report, National Security Agency, Information Assurance Directorate.

[3] Arkfeld, M. R. (2005), *Electronic Discovery and Evidence*, Law Partner Publishing, L.L.C..

[4] Carrier, B. (2005), *File System Forensic Analysis*, Addison-Wesley Professional.

[5] Cox, R. & others (1996), 'AFNI: software for analysis and visualization of functional magnetic resonance neuroimages', *Computers and Biomedical Research* **29**(3), 162-173.

[6] Manes, G. W & E. D. A. B. D. G. J. H. (2007),'A Framework for Redacting Digital Information from Electronic Devices', *IEEE Information Assurance and Security Workshop*', 56-60.

[7] Marcus, D. S, Oslen, T. R., Ramaratnam, M. and Buckner, R. L. (2007),  The Extensible Neuroimaging Archive Toolkit: an informatics platform for managing, exploring, and sharing neuroimaging data, *Neuroinformatics* **5**(1), 11-34.

[8] Russinovich, M. E. & Solomon, D. A. (2004), *Microsoft Windows Internals (4th Edition): Microsoft Windows Server 2003, Windows XP, and Windows 2000*, Microsoft Press.

[9] Silberschatz, A.; Galvin, P. B. & Gagne, G. (2004), *Operating System Concepts (7th Edition)*, Wiley.