

# Sequencing the Mouse Genome for the Oxidatively Modified Base 8-Oxo-7,8-dihydroguanine by OG-Seq

Yun Ding,<sup>1b</sup> Aaron M. Fleming,<sup>1b</sup> and Cynthia J. Burrows<sup>\*1b</sup>

Department of Chemistry, University of Utah, 315 South 1400 East, Salt Lake City, Utah 84112-0850, United States

**S** Supporting Information

**ABSTRACT:** Oxidative damage to the genome can yield the base 8-oxo-7,8-dihydroguanine (OG). In vitro studies suggested OG would preferentially form in 5'-GG-3' sequence contexts after exposure to reactive oxygen species. Herein, OG locations in the genome were studied by development of "OG-Seq" to sequence OG sites via next-generation sequencing at ~0.15-kb resolution. The results of this study found ~10 000 regions of OG enrichment in WT mouse embryonic fibroblasts and ~18 000 regions when the OG repair glycosylase Ogg1 was knocked out. Gene promoters and UTRs harbor more OG-enriched sites than expected if the sites were randomly distributed throughout the genome and correlate with reactive 5'-GG-3' sequences, a result supporting decades of in vitro studies. Sequencing of OG paves the way to address chemical and biological questions surrounding this modified DNA base, such as its role in disease-specific mutations and its epigenetic potential in gene regulation.

We report sequencing of the mouse genome for the oxidatively modified base 8-oxo-7,8-dihydroguanine (OG) by a method we are naming OG-Seq. OG stems from oxidation of the guanine (G) heterocycle by cellular reactive oxygen species (ROS).<sup>1</sup> G is favored for oxidation over the other heterocyclic DNA bases because it has the lowest redox potential.<sup>1</sup> One-electron oxidation studies of oligomers in vitro have determined that the 5' G of 5'-GG-3' sequences is preferentially oxidized; the sequence dependency in the reaction results from base stacking decreasing the redox potential at the 5' G.<sup>1,2</sup> Further, oxidations under conditions that model the reducing cellular state show increased yields of OG.<sup>1</sup> This observation is supported by extraction of genomic DNA from ROS-stressed cells where high yields of OG were identified.<sup>3</sup> A drawback to the current methods of analyzing cellular OG levels is that the DNA is digested to nucleosides followed by mass spectrometry analysis causing all sequence information to be lost. Thus, it remains unknown if the OG formed in the genome occurs randomly, or if oxidations are site or region specific. More importantly, tracking preferential locations of OG formation is critical for understanding the molecular basis of disease and cancer causing G→T transversions;<sup>4</sup> this mutation type is indirectly assigned to result from OG. Recently, evidence supporting a regulatory function for OG located in gene promoters was documented that suggests this modification may be epigenetic-like;<sup>5–8</sup> therefore, to better understand both

the mutational and epigenetic-like properties of OG, an OG sequencing approach is urgently needed.

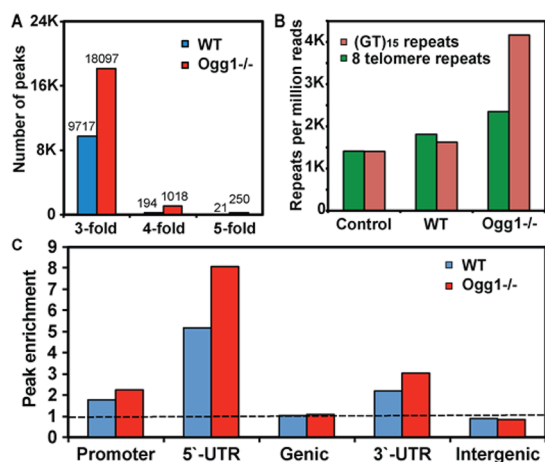
Initial attempts to sequence OG from a genome followed two different methods. The first method harnessed a PCR primer to conduct polymerase extension to elucidate if an OG-specific mutation (G→T) occurs at a given site by Sanger sequencing.<sup>7,9</sup> The second type of method utilized an OG antibody to find OG sites in the genome. If the genomic DNA is fragmented followed by OG antibody enrichment, microarray analysis allows OG sequencing at ~10-kb resolution.<sup>10</sup> When in situ immunodetection of OG was used to construct a chromosome map of OG, the resulting map was obtained at ~1000-kb resolution.<sup>11</sup> The OG antibody allowed ChIP-Seq analysis of OG in the rat genome at ~0.1-kb resolution.<sup>8</sup> These techniques all offer significant advancement of our knowledge regarding genomic OG, although each method has drawbacks. The PCR primer approach is low throughput and requires prior knowledge of critical locations to interrogate. For the antibody methods, generally the resolution is not high enough to determine precise genomic elements (i.e., promoters, UTRs, etc.) in which OG resides; additionally, as we demonstrate below, the OG antibody binding affinity is significantly impacted by DNA secondary structures. The latter issue introduces considerable bias in the data obtained. In the present work, an approach was developed to map OG in the mouse genome at ~0.15-kb resolution that utilized chemistry to label OG with biotin for affinity purification (AP) rather than immunoprecipitation (IP) for sample enrichment. The enriched samples were analyzed by next-generation sequencing (NGS) to obtain the data sets described for the entire mouse genome.

Finding the sequences of the ~1–2 OGs per 10<sup>6</sup> G's for the entire mouse genome<sup>11</sup> requires fragmenting the genome to smaller pieces (~150 mers) followed by enrichment of the OG-containing strands. To identify the best enrichment protocol, control experiments were conducted with synthetic OG-containing oligomers of known sequence and secondary structure. The controls were also spiked with duplex oligomers devoid of OG to determine the impact that the non-OG-containing strands from the genomic sample (~10<sup>4</sup>–10<sup>5</sup> more) will have on the enrichment efficiency. A commercially available OG antibody was titrated from 10 to 7000 equiv relative to OG in the contexts of ssDNA, dsDNA, and G4-DNA (Figures S1 and S2). These structures represent common contexts in which OG will reside after the genomic DNA has been fragmented and denatured by flash cooling for IP. We monitored the capture by

**Received:** December 7, 2016

**Published:** February 2, 2017





**Figure 3.** Analysis of NGS reads obtained by applying OG-Seq to WT and Ogg1<sup>-/-</sup> MEFs. (A) Peak counts from the two cell lines at 3-, 4-, and 5-fold enrichment and above. (B) Counts of (5'-GT-3')<sub>n</sub> microsatellite and telomere repeat DNA per million unmapped reads extracted from the two cell lines. (C) Enrichment in the 3-fold or more enriched peaks in each genomic element of the WT and Ogg1<sup>-/-</sup> MEF genomes. The counts plotted represent the fold enrichment observed in the genomic element relative to the expected count if the OG-enriched peaks were randomly distributed throughout the genome on the basis of each element's genomic distribution.

(Figure 2B). Starting at 3-fold enrichment, 9717 peaks in the WT-cell genomes and 18 097 from the Ogg1<sup>-/-</sup>-cell genomes were observed (Figure 3A). The nearly 2-fold increase in peaks observed in the Ogg1<sup>-/-</sup> cells is consistent with the ~2-fold increase in OG observed in the mice from which these cell lines were obtained.<sup>15</sup> At 4-fold enrichment, there were many fewer peaks overall, and the Ogg1<sup>-/-</sup>-cell genomes had 5× more peaks than the WT-cell genomes (Figure 3A). Lastly, at 5-fold enrichment, there again was a reduction in the peak count and the Ogg1<sup>-/-</sup>-cell genomes had nearly 10× more peaks of OG enrichment over the WT-cell genomes (Figure 3A). The observation that sequencing reads aligned to give peaks of enrichment supports a conclusion that OG formation is nonrandom in the genome. Moreover, the genomic regions enriched in OG, as is evident by the peaks, are regions that are more highly reactive toward oxidation than those where OG-enriched peaks were not observed. Although each set of sequencing experiments, with WT vs Ogg1<sup>-/-</sup> genomes, was only performed once, the fact that the increase in OG in the repair-compromised cells was consistent with literature reports of the overall 2-fold increase in OG lends confidence to the results.<sup>15</sup> Further experiments will be needed to verify the consistency of the sites of modification in various cell lines and at different times in the cell cycle.

Repeat DNA does not align to the reference genome and is generally discarded from NGS analysis; however, we counted the number of reads comprised of the telomere repeat sequence (5'-TTAGGG-3')<sub>n</sub> and the (5'-GT-3')<sub>n</sub> microsatellites because they have the same percent G but different sequence contexts and chromosome positions. The counts are presented as number of repeat reads per million unmapped reads obtained. We found the Ogg1<sup>-/-</sup>-cell genomic DNA contained more repeat strands with OG than the WT-cellular genomes; additionally, the (5'-GT-3')<sub>n</sub> repeat count doubled while the telomere count increased by 33% in the Ogg1<sup>-/-</sup>-cell genomes relative to the WT. We had anticipated more OG in the telomeres than the (5'-GT-3')<sub>n</sub> on

the basis of our previous experiments;<sup>18</sup> however, the total lengths of these repeats in the mouse genome are not yet well understood,<sup>19</sup> and telomeres show great variability in length between chromosomes and cells<sup>20</sup> that prevents us from making any strong conclusions from these values. Lastly, telomeres are generally exposed to more damage<sup>20</sup> likely resulting in further oxidation of OG to hydantoin products<sup>18</sup> that would be silent in the present sequencing experiment.

In the next analysis, the genomic regions with peaks of 3-fold enrichment or more were inspected based on the genomic element in which they reside (Figures 2B and 3C), including promoters, 5'-UTRs, 3'-UTRs, genic regions (exons plus introns), and intergenic regions. The data in Figure 3C were normalized with respect to the relative distribution each genomic element has in the mouse genome (Figure S6). First, between the WT- and Ogg1<sup>-/-</sup>-cell genomes, intergenic regions provided fewer peaks than would be predicted if enriched peaks were randomly distributed. This observation is consistent with intergenic regions being protected as heterochromatin to safeguard them from oxidation. Interestingly, promoters, 5'-UTRs, and 3'-UTRs provided greater relative numbers of OG-enriched peaks than expected by a random distribution of the peaks throughout the genome. Additionally, the genomes from Ogg1<sup>-/-</sup> cells gave more OG-enriched peaks (1.3–1.6×) than the WT cells in these reactive regions. The promoters and UTRs are critical for gene regulation causing them to be the most exposed for regulatory protein interactions. Consequently, the increased exposure of these genomic elements apparently results in their increased levels of G oxidation to OG. The genes found enriched with OG are provided in Tables S1 and S2. Notable examples from these tables include the oncogenes *Brca1*, *c-Kit*, *Ret*, and *Palm*. The genic regions from both cell lines provided peaks of OG enrichment that were similar in distribution as expected if the peaks were randomly distributed throughout the genome. It is not clear why genic regions were less reactive than the flanking control UTRs, but the importance of maintaining a proper coding sequence must cause the cell to guard these regions from oxidation.

A chromosomal-level analysis of the OG-enriched peaks (3-fold enrichment) was then made for the WT and Ogg1<sup>-/-</sup> genomes (Figure S7). First, the number of OG-enriched peaks per chromosome was not dependent on the chromosome length (Figure S7). For example, chromosome 10 provided the most OG-enriched peaks from the two cell lines, but this chromosome is intermediate in length, and a similar observation occurs with chromosome 15. Next, we looked to see if chromosomes that are richer in genes gave more OG-enriched peaks. Chromosomes possessing a greater density of genes possibly have greater amounts of euchromatin and may be more reactive toward oxidation. For instance, chromosomes 7 and 11 are gene rich but were found to have some of the lowest levels of OG-enriched peaks; thus, greater gene density does not yield greater G oxidation to OG.

In vitro oxidations have found the 5' G in 5'-GG-3' sequence contexts are more reactive toward oxidation.<sup>2,21</sup> Thus, the percentage of 5'-GG-3' sequences in each chromosome was determined and compared to the enriched peak counts from each chromosome (Figure S7). Chromosomes 7 and 11 have the largest percentage of global 5'-GG-3' sequences; however, these chromosomes did not show greater amounts of OG-enriched peaks. Thus, G content must not define chromosome reactivity toward oxidation. Although, when the OG-enriched peaks were inspected for their percentage of 5'-GG-3' sequences, we found

that the peaks were composed of more of these local reactive sequences than expected if the sequences were randomly selected from the genome (Figure S8). The present results provide genome-level results supporting the many in vitro studies that predicted oxidation of 5'-GG-3' sequences would dominate in vivo.<sup>1,2,21</sup> Lastly, one hypothesis for this non-random chromosome-level distribution of OG-enriched peaks may reside in the non-random spatial distribution of interphase chromosomes in the nucleus.<sup>10</sup> Interphase chromosomes bias regions for interacting with the nuclear envelope causing these regions to be preferentially exposed to more oxidants diffusing into the nucleus, while other regions are protected from oxidation because they are toward the interior of the nucleus.<sup>10</sup> The non-random spatial distribution of chromosomes in the nucleus would bias the observed regions of G oxidation.

The OG-enriched peaks were then inspected for G4's because these G-rich sequences should be more prone to oxidation on the basis of our previous studies.<sup>18</sup> The analysis found ~20% of the peaks from WT-MEF genomes and ~25% of the *Ogg1*<sup>-/-</sup>-MEF genomes possessed potential G4's. These G4 counts were greater than expected when compared to randomized samples (Figure S9). An increased level of OG enrichment in G4's was also documented by Gillespie et al.<sup>8</sup>

The present work sequenced OG in the mouse genome via OG-Seq developed in this work. At the heart of this method is the power of chemistry to label OG with BTN for STP enrichment (i.e., AP) while minimizing the structural bias that is a major limitation for IP enrichment protocols (Figure 2A). Enrichment of the fragmented and labeled duplexes by AP followed by release of the complementary strands allowed OG-Seq at 0.15-kb resolution. The method was applied to WT- and *Ogg1*<sup>-/-</sup>-MEF genomes to find regions of OG enrichment. These regions preferentially reside in promoter and UTR regulatory regions flanking protein-coding sequences (Figure 3C). These genomic elements exist in euchromatin regions that are less protected resulting in greater oxidation of G to OG. Additionally, the OG-enriched peaks harbored more 5'-GG-3' reactive sequences and G4's than expected by random chance (Figure S7). These findings support many decades of in vitro studies aimed at understanding oxidative damage to cellular DNA.<sup>1,2,21</sup>

We further recognize that sequencing of the biotinylated strands would introduce a characteristic mutation at the Sp-BTN nucleotide after polymerase bypass (Figures 1B and 2A). During NGS library preparation, a PCR step is performed. The polymerase bypass of the Sp-BTN adduct would most likely yield a characteristic G→T and G→C mutation signature if these sites are bypassed in the same way as the spiroiminodihydroantoin core structure.<sup>22</sup> Mining the sequence reads from an NGS experiment for these characteristic mutations would allow single-nucleotide resolution of the OG locations. Critical to the ability to successfully achieve this goal is sequencing at high depth (>30×) to perform proper statistical analysis of the results. Future application of this approach will greatly increase our knowledge surrounding the chemical biology of OG with respect to its disease-causing mutation potential<sup>4</sup> and point to possible sites where OG plays an epigenetic-like regulatory role.<sup>5-7</sup>

## ■ ASSOCIATED CONTENT

### 📄 Supporting Information

The Supporting Information is available free of charge on the ACS Publications website at DOI: 10.1021/jacs.6b12604.

Detailed methods, HPLC chromatograms for biotinylation of each sequence context, and NGS data analysis, including Figures S1–S8 and Tables S1 and S2 (PDF)

## ■ AUTHOR INFORMATION

### Corresponding Author

\*burrows@chem.utah.edu

### ORCID

Yun Ding: 0000-0002-6624-0934

Aaron M. Fleming: 0000-0002-2000-0310

Cynthia J. Burrows: 0000-0001-7253-8529

### Notes

The authors declare no competing financial interest.

## ■ ACKNOWLEDGMENTS

The MEF cell lines were provided by Dr. Tomas Lindahl (Imperial Cancer Research Fund, UK). We thank the University of Utah core facilities for assistance. This work was partially supported by a National Cancer Institute grant (R01 CA090689) and a seed grant from the Nuclear Control Program at Huntsman Cancer Institute at the University of Utah (P30 CA042014).

## ■ REFERENCES

- (1) Fleming, A. M.; Burrows, C. J. *Free Radical Biol. Med.* **2016**, DOI: 10.1016/j.freeradbiomed.2016.11.030.
- (2) Genereux, J. C.; Barton, J. K. *Chem. Rev.* **2010**, *110*, 1642.
- (3) Cadet, J.; Douki, T.; Ravanat, J.-L. *Mutat. Res., Fundam. Mol. Mech. Mutagen.* **2011**, *711*, 3.
- (4) Roberts, S. A.; Gordenin, D. A. *Nat. Rev. Cancer* **2014**, *14*, 786.
- (5) Fleming, A. M.; Ding, Y.; Burrows, C. J. *Proc. Natl. Acad. Sci. U. S. A.* **2016**, DOI: 10.1073/pnas.1619809114.
- (6) Pan, L.; Zhu, B.; Hao, W.; Zeng, X.; Vlahopoulos, S. A.; Hazra, T. K.; Hegde, M. L.; Radak, Z.; Bacsı, A.; Brasier, A. R.; Ba, X.; Boldogh, I. J. *Biol. Chem.* **2016**, *291*, 25553.
- (7) Park, J.; Park, J. W.; Oh, H.; Maria, F. S.; Kang, J.; Tian, X. *PLoS One* **2016**, *11*, e0155792.
- (8) Pastukh, V.; Roberts, J. T.; Clark, D. W.; Bardwell, G. C.; Patel, M.; Al-Mehdi, A. B.; Borchert, G. M.; Gillespie, M. N. *Am. J. Physiol. Lung Cell Mol. Physiol.* **2015**, *309*, L1367.
- (9) Nachtergaeel, A.; Belayew, A.; Duez, P. *DNA Repair* **2014**, *22*, 147.
- (10) Yoshihara, M.; Jiang, L.; Akatsuka, S.; Suyama, M.; Toyokuni, S. *DNA Res.* **2014**, *21*, 603.
- (11) Ohno, M.; Miura, T.; Furuichi, M.; Tominaga, Y.; Tsuchimoto, D.; Sakumi, K.; Nakabeppu, Y. *Genome Res.* **2006**, *16*, 567.
- (12) Hosford, M. E.; Muller, J. G.; Burrows, C. J. *J. Am. Chem. Soc.* **2004**, *126*, 9540.
- (13) Xue, L.; Greenberg, M. M. *J. Am. Chem. Soc.* **2007**, *129*, 7010.
- (14) Bajacan, J. E. V.; Hong, I. S.; Penning, T. W.; Greenberg, M. M. *Chem. Res. Toxicol.* **2014**, *27*, 1227.
- (15) Klungland, A.; Rosewell, I.; Hollenbach, S.; Larsen, E.; Daly, G.; Epe, B.; Seeberg, E.; Lindahl, T.; Barnes, D. E. *Proc. Natl. Acad. Sci. U. S. A.* **1999**, *96*, 13300.
- (16) David, S. S.; O'Shea, V. L.; Kundu, S. *Nature* **2007**, *447*, 941.
- (17) Fleming, A. M.; Armentrout, E. I.; Zhu, J.; Muller, J. G.; Burrows, C. J. *J. Org. Chem.* **2015**, *80*, 711.
- (18) Fleming, A. M.; Burrows, C. J. *Chem. Res. Toxicol.* **2013**, *26*, 593.
- (19) Ellegren, H. *Nat. Rev. Genet.* **2004**, *5*, 435.
- (20) O'Sullivan, R. J.; Karlseder, J. *Nat. Rev. Mol. Cell Biol.* **2010**, *11*, 171.
- (21) Liu, Y.; Liu, Z.; Geacintov, N. E.; Shafirovich, V. *Phys. Chem. Chem. Phys.* **2012**, *14*, 7400.
- (22) Kornysheva, O.; Burrows, C. J. *Biochemistry* **2003**, *42*, 13008.