

# Sources of NGS Error and What to Do About Them

# Standard Illumina Library Prep

1. Harvest gDNA
2. Shear gDNA
3. Blunt-end-repair DNA
4. dA-tail DNA
5. Ligate Adapters to DNA
6. PCR Amplify Library
7. Sequence

# DNA Damage source

Source of DNA	Potential Damage	Comments	References
Ancient DNA	abasic sites, deaminated cytosine, oxidized bases, fragmentation, nicks	Cytosine deamination has been reported to be the most prevalent cause of sequencing artifacts in ancient DNA.	Gilbert, M.T. et al. (2007) <i>Nuc. Acid Res.</i> , 35, 1–10. Hofreiter, M. et al. (2001) <i>Nuc. Acid Res.</i> , 29, 4793.
Environmental DNA	fragmentation, nicks (plasmid or genomic)	Nicks and fragmentation can increase the formation of artifactual chimeric genes during amplification.	Qiu, X. et al. (2001) <i>Appl. Envir. Microbiol.</i> , 67, 880.
Source of Damage			
Exposure to Ionizing Radiation	abasic sites, oxidized bases, fragmentation, nicks	Ionizing radiation is used to sterilize samples.	Sutherland, B.M. et al. (2000) <i>Biochemistry</i> , 39, 8026.
Exposure to Heat	fragmentation, nicks, abasic sites, oxidized bases, deaminated cytosine, cyclopurine lesions	Heating DNA accelerates the hydrolytic and oxidative reactions in aqueous solutions.	Bruskov, V.I. (2002) <i>Nuc. Acids Res.</i> , 30, 1354.
Phenol/Chloroform Extraction	oxidized bases	Guanine is more sensitive to oxidation than the other bases and forms 8-oxo-G. 8-oxo-G can base pair with A making this damage potentially mutagenic.	Finnegan, M.T. (1995) <i>Biochem. Soc. Trans.</i> , 23, 403S.
Exposure to Light (UV)	thymine dimers, (cyclobutane pyrimidine dimers) pyrimidine (6–4) photo products	UV trans-illumination to visualize DNA causes thymine dimer formation.	Cadet, J. et al. (2005) <i>Mutat. Res.</i> , 571, 3–17. Pfeifer, G.P. et al. (2005) <i>Mutat. Res.</i> , 571, 19–31.
Mechanical Shearing	fragmentation, nicks	Normal DNA manipulations such as pipetting or mixing can shear or nick DNA.	
Dessication	fragmentation, nicks, oxidized bases		Mandrioli, M. et al. (2006) <i>Entomol. Exp. App.</i> , 120, 239.
Storage in Aqueous Solution	abasic sites, oxidized bases, deaminated cytosine, nicks, fragmentation	Long term storage in aqueous solution causes the accumulation of DNA damage.	Lindahl, T. et al. (1972) <i>Biochemistry</i> , 11, 3610 and 3618.
Exposure to Formalin	DNA-DNA crosslinks, DNA-protein crosslinks	Formaldehyde solution that has not been properly buffered becomes acidic, increasing abasic site formation.	Workshop on recovering DNA from formalin preserved biological samples. (2006) The National Academies Press.

From NEB expressions Spring 2007, vol 2.1  
By Thomas C. Evans, Jr., New England Biolabs, Inc.

[Link](#)

# Sources of Errors in Illumina Library Prep

1. Harvest gDNA
  2. Shear gDNA
  3. Blunt-end-repair DNA
  4. dA-tail DNA
  5. Ligate Adapters to DNA
  6. PCR Amplify Library
  7. Sequence
1. Outgrowth, Storage
  2. DNA Oxidation
    1. [Costello et al 2013 NAR](#)
  3. T4 DNA Pol est.  $1 \times 10^{-5}$ 
    1. [Keohavong, Thilly 1989 PNAS](#)
  4. Interactions with ligation?
  5. ~11% 5' anti-T pro A/G bias
    1. [Seguin-Orlando et al 2013 PLOS ONE](#)
  6. Phusion  $4.2 \times 10^{-7}$ 
    1. [Li et al 2006 Nature Methods](#)
  7. Sequence specific, PCR
    1. [Nakamura et al 2011 NAR](#)

# Final Results

- Overall error rate estimated between 0.1 and 1% per base ([Lou et al 2013 PNAS](#)).
- That's 1 error per 100-1000 bases sequenced, or typically at least 1 error per 5 paired reads.
- 2 billion 100bp PE reads / run means between 400,000,000 and 4,000,000,000 errors per run.
- Minimum detection limit is between 0.1 and 1%

**IF EVERYTHING WE DO GENERATES  
ERRORS WHAT CAN WE DO?**

# Some Suggestions to Minimize Errors

1. Minimize sample handling after biological relevance.
2. Minimal PCR whenever possible.
3. Pay attention to directionality of reads supporting variant.
4. Make use of quality score information.
5. Use sequence specific error profiles to eliminate false positives.
  1. [Meacham et al 2011 BMC Bioinformatics](#)
6. Leverage other biological knowledge whenever possible (ie timecourse data).

# A Note On False Negatives

- CNV between repeat elements can be virtually invisible (particularly at low levels).



# Practical Limitations

- If planned sequencing coverage is less than  $\sim 100$ , most not important (except sequence specific). Always assume something seen once is not real.
- If looking for phenotypes, driver mutations in cancer, or other disease associated mutations, causal mutations not likely to be rare.
- Massaging standard illumina data is likely to be less effective than better experimental planning and design using alternative library preparation methods.