

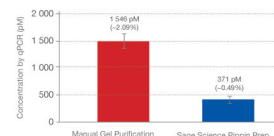
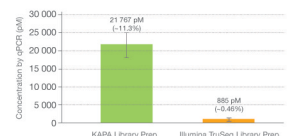
GAVIN J. RUSH<sup>1</sup>, ERIC VAN DER WALT<sup>1</sup>, JACOB KITZMAN<sup>2</sup>, LIESL NOACH<sup>1</sup>, ZAYED ALBERTYN<sup>3</sup>, COLIN HERCUS<sup>3</sup>, CHARLIE LEE<sup>2</sup>, JAY SHENDURE<sup>2</sup>, JOHN F. FOSKETT III<sup>1</sup>, PAUL J. MCEWAN<sup>1</sup>  
<sup>1</sup> KAPABIOSYSTEMS, 600 WEST CUMMINGS PARK, SUITE 2250, WOBURN, MA 01801 | <sup>2</sup> UNIVERSITY OF WASHINGTON, DEPARTMENT OF GENOME SCIENCES, FOGG BUILDING S-250, BOX 355065, 3720 15TH AVE NE, SEATTLE WA 98195 | <sup>3</sup> NOVOGRAFT, C-23A-5, TWO SQUARE, SECTION 19, 46300 PETALING JAYA, SELANGOR, MALAYSIA

## SAMPLE PREPARATION

### INTRODUCTION

Dramatic improvements in commercial Next Generation Sequencing (NGS) platforms have resulted in spectacular reductions in the cost-per-base of DNA sequencing. Until recently, the primary focus for innovation has been on improvements to the core sequencing technologies, with optimization of sample preparation playing a secondary role. The exponential gains in sequencing capacity have simultaneously led to higher sample throughput, placing increasing emphasis on the importance of improved library construction protocols for multiplexed sample sequencing. While major commercial NGS systems all require the construction of similar libraries via analogous workflows, some protocols and/or reagents offer significant advantages over others, and end-users must choose among numerous alternative methods and reagents for sample preparation. We re-sequenced the *Staphylococcus aureus*, *Escherichia coli*, and *Mycobacterium tuberculosis* genomes to compare the standard Illumina TruSeq sample preparation reagents and workflow with a number of innovative improvements including: alternative library preparation reagents and protocols; automated fragment size selection; real-time library amplification; amplification-free sequencing; and accurate qPCR library quantification for sample pooling and multiplexed sequencing.

SPECIES	GENBANK ID	GENOME SIZE (bp)	% GC
<i>S. aureus</i> subsp. <i>aureus</i> TW20	FN43396	3,043,210	33%
<i>E. coli</i> DH10B	NC_010473	4,696,137	51%
<i>M. tuberculosis</i> H37Rv	AL23456	4,411,532	65%



#### Illumina TruSeq Low-Throughput (LT) Protocol

1 µg of sheared gDNA per library

END REPAIR	CLEAN-UP (AMPURE XP BEADS)	A-TAILING	LIGATE ADAPTORS	NO CLEAN-UP	LIGATE ADAPTORS
1 µg sheared gDNA in PBS: 50 µL End Repair Control DNA: 40 µL End Repair Buffer (KAPA): 40 µL Total: 130 µL Incubation: 30 min @ 30 °C	Illumina TruSeq Low-Throughput protocol	End Reagent gDNA: 15 µL A-Tailing Control DNA (Illumina): 15.8 µL A-Tailing Mix (Illumina): 35 µL Total: 65.8 µL Incubation: 30 min @ 30 °C	A-Tail and End Reagent gDNA: 30 µL Ligase Control DNA (Illumina): 2.5 µL DNA Ligase Mix (Illumina): 12.5 µL DNA Adaptor Index 2, 4, 5, 6 or 7 (Illumina): 2.5 µL Total: 50 µL Incubation: 15 min @ 30 °C Add 5 µL Stop Ligase Mix (Illumina)	No clean-up as per Illumina TruSeq LT protocol	A-Tail and End Reagent gDNA: 30 µL Ligase Control DNA (Illumina): 2.5 µL DNA Ligase Mix (Illumina): 12.5 µL DNA Adaptor Index 2, 4, 5, 6 or 7 (Illumina): 2.5 µL Total: 50 µL Incubation: 15 min @ 30 °C No Stop - diversity to AMPure XP clean-up

#### KAPA Library Prep Protocol

1 µg of sheared gDNA per library

END REPAIR	CLEAN-UP (AMPURE XP BEADS)	A-TAILING	LIGATE ADAPTORS	NO CLEAN-UP	LIGATE ADAPTORS
1 µg sheared gDNA in PBS: 50 µL End Repair Control DNA: 40 µL End Repair Buffer (KAPA): 40 µL Total: 130 µL Incubation: 30 min @ 30 °C	Illumina TruSeq Low-Throughput protocol	End Reagent gDNA: 15 µL A-Tailing Control DNA (Illumina): 15.8 µL A-Tailing Mix (Illumina): 35 µL Total: 65.8 µL Incubation: 30 min @ 30 °C	A-Tail and End Reagent gDNA: 30 µL Ligase Control DNA (Illumina): 2.5 µL DNA Ligase Mix (Illumina): 12.5 µL DNA Adaptor Index 2, 4, 5, 6 or 7 (Illumina): 2.5 µL Total: 50 µL Incubation: 15 min @ 30 °C Add 5 µL Stop Ligase Mix (Illumina)	No clean-up as per Illumina TruSeq LT protocol	A-Tail and End Reagent gDNA: 30 µL Ligase Control DNA (Illumina): 2.5 µL DNA Ligase Mix (Illumina): 12.5 µL DNA Adaptor Index 2, 4, 5, 6 or 7 (Illumina): 2.5 µL Total: 50 µL Incubation: 15 min @ 30 °C No Stop - diversity to AMPure XP clean-up

#### Yield of adaptor-ligated library molecules before size selection.

Sheared genomic DNA was prepared in bulk by nebulization and identical starting material was used for each library (1 µg of gDNA from *S. aureus*, *E. coli*, or *M. tuberculosis*). Libraries were constructed using reagents and procedures (see detailed methods on left-hand side) from KAPABiosystems (green; 3 libraries), or using the Illumina TruSeq DNA Sample Prep Kit and Low Throughput Protocol (orange; 15 libraries). Libraries were quantified by qPCR before size selection using the KAPABiosystems Library Quantification Kit according to the recommended protocol. The estimated percentage of starting material that was converted to useful, adaptor-ligated (PCR-amplifiable) library molecules is provided.

The KAPABiosystems library construction reagents and protocol produced ~25-fold more adaptor-ligated library fragments from the same amount of starting material than did the standard Illumina TruSeq DNA Sample Prep Kit. Important differences between the two library construction procedures assessed here include: Combined enzyme and reaction buffer "master mix" formulations (TruSeq vs. separate enzyme and reaction buffer (KAPA)); and the absence of an Ampure XP Bead clean-up between A-tailing and Adaptor Ligation in the TruSeq Protocol.

#### MANUAL GEL EXTRACTION

- Requires manual gel excision and downstream extraction of the sample from the agarose gel.
- Actual DNA fragment length distribution is broader and less reproducible.
- Manual gel electrophoresis requires longer run time and more manual steps.
- Contamination of adaptor-dimers and other low molecular weight artifacts leads to lower data quality.
- Potential for cross-contamination amongst samples on a single gel.

#### AUTOMATED SIZE SELECTION

- Automated elution of sample material into buffer downstream with cleanable workflows (PCR amplification, flow cell amplification).
- Selects tight, accurate, and reproducible DNA fragment size ranges.
- Significant time and labor savings from quicker run times and automated elution.
- Reduction of low molecular weight DNA contamination reduces library artifacts.
- Reduced potential for cross-contamination amongst samples.

**MANUAL GEL PURIFICATION**

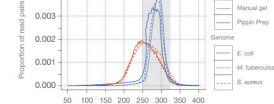
Adaptor gDNA:	Adaptor gDNA:	Adaptor gDNA:
20 µL	20 µL	20 µL
10 µL	10 µL	10 µL
40 µL	40 µL	40 µL
Total:	Total:	Total:
70 µL	70 µL	70 µL

2% agarose gel in 1x TBE  
 Size range collected: 270 bp - 450 bp  
 Run time: 15 min

#### Yield of adaptor-ligated library molecules after size selection.

Starting with 1 µg of either *S. aureus*, *E. coli*, or *M. tuberculosis* gDNA, we constructed duplicate libraries using reagents and procedures from the KAPA DNA Library Prep Kit (see detailed methods on left-hand side). Duplicate libraries were pooled, and the three resulting libraries (*S. aureus*, *E. coli*, or *M. tuberculosis*) were each split equally before size selection either by manual gel purification (Qiagen QIAquick Gel Extraction Kit) or by Pippin Prep automated gel purification (Sage Science). Libraries were quantified by qPCR using the KAPABiosystems Library Quantification Kit according to the recommended protocol. The estimated percentage of starting material converted to useful, adaptor-ligated (PCR-amplifiable) library molecules is provided.

Size selection by manual gel purification yielded ~4-fold more DNA than size selection by automated gel purification using the Sage Science Pippin Prep apparatus. To a large extent this may be accounted for by the fact that the Pippin Prep method yielded a much narrower size range of library fragments (see below), despite careful manual gel purification targeting an equivalent size range of library molecules.



#### Size distribution of library molecules after size selection by manual agarose gel purification or automated gel purification with Sage Science Pippin Prep.

Starting with 1 µg of either *S. aureus*, *E. coli*, or *M. tuberculosis* sheared gDNA, we constructed duplicate libraries using reagents and procedures from the KAPA DNA Library Prep Kit (see detailed methods on left-hand side). Duplicate libraries were pooled, and the three resulting libraries (*S. aureus*, *E. coli*, or *M. tuberculosis*) were each split equally before size selection either by manual gel excision and purification (Qiagen QIAquick Gel Extraction Kit) or by Pippin Prep automated gel purification (Sage Science). We set the Pippin Prep instrument to collect library molecules in the 370 bp - 450 bp range including adaptors, and we attempted to select the same range via manual agarose gel purification. TruSeq adaptors are 121 bp in total, suggesting targeted insert sizes were ~250 bp - 330 bp (broken line). Paired reads (2 x 75 bp) were used to determine actual insert sizes.

## AMPLIFICATION

### RATIONALE FOR REAL-TIME HIGH FIDELITY AMPLIFICATION OF NEXT-GENERATION DNA SEQUENCING LIBRARIES

High fidelity PCR is used to selectively enrich library fragments carrying appropriate adaptor sequences and to amplify the amount of DNA prior to sequencing. During PCR enrichment of libraries, a subset of library molecules are amplified with reduced efficiencies, introducing bias and resulting in uneven sequence coverage. GC content is known to be an important factor in NGS library amplification bias, and different PCR enzymes and buffer formulations are likely to show individual strengths and weaknesses in this regard. Furthermore, such biases - along with other artifacts such as PCR-induced errors, adaptor dimers, PCR duplicates, and chimeras - are exacerbated by over-amplification, while under-amplification results in insufficient yields. Inherent uncertainty in the outcome of end-point PCR often demands downstream validation of library quality by electrophoresis.

Here we present a high fidelity, real-time PCR method for rapid and convenient enrichment and amplification of libraries. The benefits of this approach include: 1) automatable workflows, 2) built-in quality metrics for each enriched library, eliminating expensive and time-consuming post-enrichment gel electrophoresis, 3) precise control over the number of PCR cycles required for optimal amplification, 4) a quality control metric for identifying inconsistencies in library preparation, and 5) seamless integration with KAPA qPCR Library Quantification kits.

#### TruSeq Low-Throughput Library Amplification Workflow

Approximate duration: 5 min

**SIZE SELECTION**

Size selection using either:  
 1. Manual gel excision, purification (Qiagen QIAquick Gel Extraction Kit, Qiagen). Recover 20 µL library DNA.  
 2. Automated DNA size selection and collection system (Pippin Prep, Sage Science). Recover 20 µL library DNA.

**PROGRAM THERMOCYCLER**

98° for 30 s  
 18 cycles of 98° for 15 s  
 67° for 30 s  
 72° for 30 s

**SET UP AND INITIATE SMALL-SCALE qPCR**

No gel electrophoresis in protocol.  
 Note: This is a targeted cycle number to determine the optimal cycle number for the subsequent preparative PCR reaction.

**REMOVE SMALL-SCALE PCR REACTIONS AT DESIRED PCR CYCLES**

Place the thermocycler and remove a small scale PCR tube after each predetermined cycle and place on ice before initiating cycling.

**VALIDATE IDEAL CYCLE NUMBER BY PERFORMING GEL ELECTROPHORESIS**

Run an aliquot of each small scale sample on agarose gel or Bioanalyzer strip to determine the optimal cycle number to be used for preparative gel electrophoresis.

**PROGRAM THERMOCYCLER FOR PREPARATIVE PCR**

98° for 30 s  
 18 cycles of 98° for 15 s  
 67° for 30 s  
 72° for 30 s

**SET UP AND INITIATE PREPARATIVE PCR**

1 x 50 µL qPCR reaction per library  
 TruSeq PCR Mix: 25 µL  
 TruSeq PCR Primer Mix: 5 µL  
 Size selected library: 20 µL

**CLEAN UP PCR REACTION**

Use AMPure XP Beads and elute in 30 µL Resuspension buffer, according to the TruSeq protocol.

**QUANTIFY AMPLIFIED LIBRARY USING ILLUMINA QUANTIFICATION PROTOCOL GUIDE**

#### Multiple Sample Workflow using KAPA High Fidelity qPCR Kit for TruSeq Library Amplification

Approximate duration: 5 min

**SIZE SELECTION**

Size selection using either:  
 1. Manual gel excision, purification (Qiagen QIAquick Gel Extraction Kit, Qiagen). Recover 20 µL library DNA.  
 2. Automated DNA size selection and collection system (Pippin Prep, Sage Science). Recover 20 µL library DNA.

**PROGRAM REAL-TIME THERMOCYCLER**

98° for 45 s  
 32 cycles of 98° for 15 s  
 67° for 30 s  
 72° for 30 s  
 Data Acquisition: 72° for 5 s\*

**SET UP AND INITIATE PREPARATIVE HIGH FIDELITY qPCR**

1 x 50 µL qPCR reaction per library  
 KAPA High Fidelity Mix: 25 µL  
 TruSeq PCR Primer Mix: 5 µL  
 Size selected library: 20 µL

**TERMINATE HIGH FIDELITY qPCR**

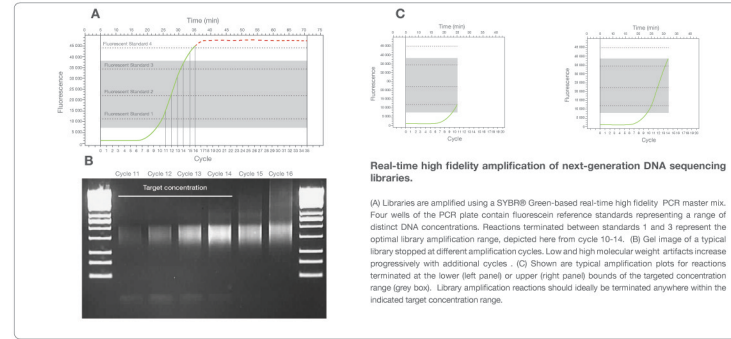
Terminate qPCR reaction when the linear amplification plots of samples for selected fluorescent standards (1-10) within targeted concentration range. This is a 4 cycle termination window and enables single step amplification of libraries with up to a 16-fold difference in initial concentration.

**CLEAN UP qPCR REACTION**

Use AMPure XP Beads and elute in 30 µL Resuspension buffer, according to the TruSeq protocol.

**QUANTIFY AMPLIFIED LIBRARY USING KAPA LIBRARY QUANTIFICATION KIT**

\*Note: This is a second step at 72° enables termination of the qPCR reaction upon the desired fluorescence intensity has been achieved before the cycle of amplification is initiated.



#### Example of multiplexed real-time high fidelity amplification.

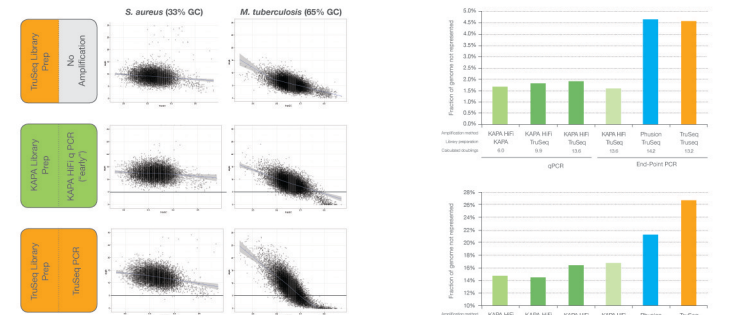
20 libraries, spanning a ~64-fold concentration range (8 cycles), were simultaneously amplified and terminated after 14 cycles. 14 of the 20 libraries fall within the targeted amplification range. The remaining 6 libraries could either be used as is, noting that they may be outside the optimal concentration range, or they could be re-amplified individually or in high- or low-concentration groups.

#### DISADVANTAGES ASSOCIATED WITH THE EXISTING TRUSeq WORKFLOW

- Wastage of library material required for cycle number optimization leads to loss of library diversity.
- Inconsistencies in reaction volume scale-up can cause variable results.
- Longer protocol involves more time and labor.
- Gel electrophoresis steps are not amenable to automation.

#### ADVANTAGES OF HIGH FIDELITY REAL-TIME PCR

- Built-in real-time quality metrics (concentration range) for each amplified DNA library.
- Real-time PCR is amenable to automation.
- Precise control over PCR cycle number required for optimal amplification.
- Seamless integration with qPCR-based library quantification.
- KAPA High Fidelity DNA Polymerase is less prone to amplification bias due to high- or low-GC content.



#### Effect of GC content on coverage depth for libraries amplified using KAPA High Fidelity qPCR Master Mix or Illumina TruSeq PCR Master Mix.

Indexed libraries were prepared from identical sheared *S. aureus* (33% GC; left panels) and *M. tuberculosis* (65% GC; right panels) gDNA using either the KAPA DNA Library Prep Kit, or the Illumina TruSeq DNA Sample Prep Kit, and then amplified using the indicated PCR reagents before paired-end sequencing (2 x 75 bp). After filtering and aligning read pairs to reference sequences, 250 000 read pairs were randomly sampled for each genome, and scatter plots of mean sequence coverage depth vs. GC content were generated by analyzing 250 bp windows.

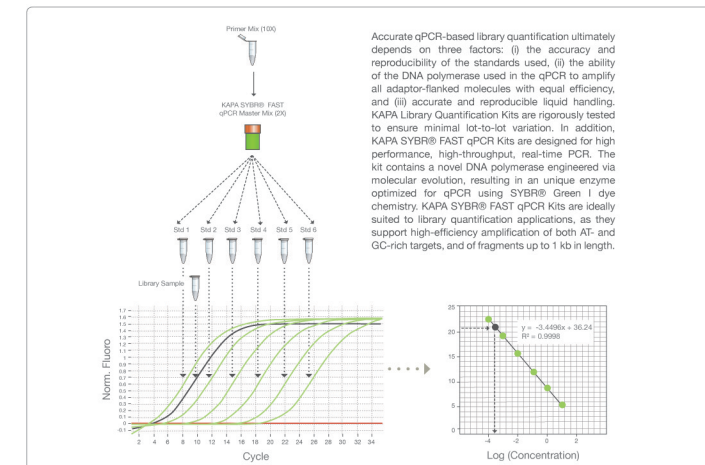
For the AT-rich *S. aureus* genome (left panel), none of the samples showed gross amplification bias compared to the unamplified control sample. GC-rich *M. tuberculosis* sequences (right panel) in the library constructed and amplified using Illumina TruSeq reagents are under-represented in the sequencing data. In contrast, the library prepared using KAPABiosystems reagents yielded coverage across the range of GC-content that is almost indistinguishable from that of the unamplified control, indicating that amplification with KAPA High Fidelity qPCR Master Mix introduced minimal additional GC-dependent coverage bias.

## LIBRARY QUANTIFICATION

### KAPA LIBRARY QUANTIFICATION KIT FOR OPTIMAL SAMPLE MULTIPLEXING

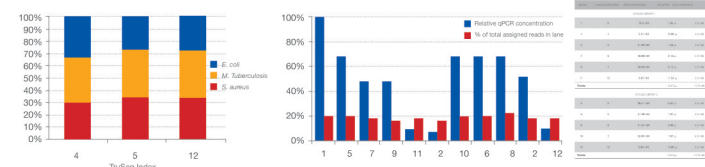
Accurate quantification of the number of amplifiable molecules in a library is critical to the outcome of sequencing results on the Illumina Genome Analyzer next-generation sequencing platform. Overestimation of library concentration results in lower cluster density after bridge PCR. Underestimation of library concentration results in too many clusters on the flow cell, which can lead to poor cluster resolution. Both scenarios result in suboptimal sequencing capacity. qPCR is widely regarded as the gold standard for accurate quantification of DNA libraries as it is the only technique capable of measuring the number of amplifiable molecules. The broad dynamic range of qPCR also enables accurate quantification of extremely dilute libraries.

- Here we demonstrate the use of the system to:
- Generate Illumina TruSeq indexed libraries comprising equimolar concentrations of 3 bacterial genomes spanning a range of average GC contents (*S. aureus*, 33% GC; *E. coli*, 51% GC; and *M. tuberculosis*, 65% GC).
  - Pool multiple Illumina TruSeq indexed libraries in equimolar concentrations for equal representation in sequence data.



#### KAPA Library Quantification Kit Workflow.

KAPA Library Quantification Kits comprise DNA Standards (six 10-fold dilutions) and 10X Primer Premix, paired with KAPA SYBR FAST qPCR Kits to accurately quantify the number of amplifiable molecules in an Illumina GA library. The 452 bp KAPA Illumina GA DNA Standard consists of a linear DNA fragment flanked by qPCR primer binding sites. Quantification is achieved by inference from a standard curve generated using the six DNA Standards.



#### qPCR quantification enables equal representation of pooled indexed libraries.

Eleven indexed Illumina TruSeq libraries were quantified by qPCR using the KAPA Library Quantification Kit according to the recommended protocol, and then combined to achieve equal final concentrations in two separate pools for multiplexed sequencing on different flow-cell lanes. The eleven libraries ranged ~11-fold in concentration from 0.67 pM to 7.65 pM, while representation of each index varied between 90% and 127% of expected assigned reads per lane.

#### qPCR quantification enables equal representation of pooled indexed libraries.

For each index (TruSeq 4, 5, and 12) we constructed three separate libraries (*S. aureus*, 33% GC; *E. coli*, 51% GC; and *M. tuberculosis*, 65% GC). Individual libraries were quantified using the KAPA Library Quantification Kit and for each index the libraries were pooled to achieve equimolar representation for each genome. The results indicate that quantification is reliable for samples with a wide range of GC content.