

1 **Detection, discrimination and discovery of a new *Tobacco streak virus* strain.**

2 **M. Dutta^a, A. Ali^b, U. Melcher^{a,c}**

3 **^a Department of Biochemistry and Molecular Biology, Oklahoma State University,**

4 **Stillwater, Oklahoma, 74078-3035 USA**

5 **^b Department of Biological Science, The University of Tulsa, Oklahoma, 74104,**

6 **USA.**

7 ^cCorresponding author: Ulrich Melcher

8 Tel.: 405-744-6210

9 FAX: 405-744-7799

10 ulrich.melcher@okstate.edu

11 **Department of Biochemistry and Molecular Biology, Oklahoma State University,**

12 **246 NRC, Stillwater, Oklahoma, 74078-3035 USA**

13

1 **Abstract:**

2 Soybean plants that exhibited symptoms of virus infection were sampled from different
3 counties of Oklahoma. These plants were tested serologically for 15 viruses known to
4 infect soybean plants. Fifty-seven samples that exhibited typical virus-like symptoms did
5 not test positive for any of the 15 viruses used in a dot-immunobinding assay (DIBA).
6 Four samples were pooled and used for next generation sequencing using the 454- Roche
7 protocol. Sequence and phylogenetic analysis of the sequences obtained revealed
8 infection with a distinct strain of *Tobacco streak virus* (TSV). TSV was one of the 15
9 viruses initially tested for using DIBA and had tested negative. TSV belongs to the genus
10 *Ilarvirus* and has been reported as a causal agent of bud blight in soybean crops in Brazil
11 and the United States. Out of 10 reported primer pairs for TSV reverse transcription-
12 polymerase chain reaction (RT-PCR), only two had the potential, based on sequence
13 similarity, to amplify part of the genome of the distinct strain of TSV found in Oklahoma
14 and only one was actually able to amplify the region. In this study, a new primer pair,
15 specific to all known TSV and capable of amplifying the Oklahoma strain (TSV-OK),
16 was designed from a highly conserved region of coat protein (CP) sequences and end-
17 point PCR and quantitative RT-PCR detection methods were developed and their
18 sensitivity assayed. This is the first report of primers designed from this specific highly
19 conserved region in the CP of TSV for detection of TSV. Twenty-three of the 57 DIBA
20 soybean samples that initially tested negative were retested with the new specific end-
21 point PCR method and found positive for TSV infection.

22

1 **Keywords:** Microbial forensic methods; 454 sequencing; RT-qPCR; SYBR green

2 **Footnote:**

3 **Abbreviations:** *Tobacco streak virus* (TSV), *Soybean mosaic virus* (SMV), *Bean pod*
4 *mottle virus* (BPMV), dot-immunobinding assay (DIBA), *Tobacco ringspot virus*
5 (TRSV), *Tomato ringspot virus* (TomRSV), movement protein (MP), coat protein (CP),
6 *Strawberry necrotic shock virus* (SNSV), *Parietaria mottle virus* (PMoV), *Blackberry*
7 *chlorotic ringspot virus* (BCRV), electron microscopy (EM), Next-generation sequencing
8 (NGS), RNA-dependent RNA polymerase (RdRp), untranslated region (UTR)

1. Introduction:

Study of diseases that may affect soybean is important to ensure good commercial crop yields. Mosaic, mottling, rugosity of leaves, mottling of seed coats, discoloration of seeds of soybean can be caused by strains of several viruses [*Tobacco streak virus* (TSV), *Soybean mosaic virus* (SMV), and *Bean pod mottle virus* (BPMV)] [1-3]. During a survey of soybean fields in Oklahoma a number of soybean plant samples were tested serologically by dot-immunobinding assay (DIBA) against the antisera of 15 different viruses including TSV (Ali et al. 2014 In preparation). More than 50% of the collected samples were negative to all 15 tested viruses. In this work, some of the previously collected DIBA negative samples (Ali et al. 2014, in preparation) were tested by NGS, which identified a novel strain of TSV that infects soybean plants in Oklahoma. TSV has been reported to cause bud blight in Brazil and the United States [4]. Bud blight causes pod necrosis and dieback of terminal portions of stems of a plant and has been observed in soybean in the past and in recent years. In the U.S., TSV was previously isolated from plants expressing symptoms of bud blight in commercial soybean fields in Wisconsin and Oklahoma [4]. Bud blight has several additional known causal agents, including *Tobacco ringspot virus* (TRSV), *Tomato ringspot virus* (TomRSV), and SMV [4]. Almost all soybean varieties tested have been susceptible [5,3]. Not much is known about the impact of TSV on soybean health and productivity and hence TSV infection merits further investigation. [4]

Viruses within genus *Ilarvirus* family *Bromoviridae*, of which TSV is the type member, have positive-sense, single-stranded RNA tripartite genomes (RNA 1 to 3). RNA 1 and

1 RNA 2 encode for the viral replication protein and RNA-dependent RNA polymerase
2 (RdRp), respectively [6]. RNA 3 encodes the movement protein (MP) and a subgenomic
3 RNA 4 that is translated into the coat protein (CP) [6].

4 Currently, there is no specified level of sequence similarity to distinguish among species,
5 strains or subgroups of ilarvirus members [7]. Serological relatedness has been used to
6 assign species to subgroups, and serological differences, host specificity or geographical
7 location have been used to distinguish different species belonging to the same genus
8 [8,9]. RNA 2 of members of subgroups 1 and 2 have a unique second open reading
9 frame absent in other ilarvirus subgroups [6]. The type isolate of TSV was originally
10 isolated from white clover (TSV-WC). Species within subgroup 1 include TSV,
11 *Strawberry necrotic shock virus* (SNSV), *Parietaria mottle virus* (PMoV), *Blackberry*
12 *chlorotic ringspot virus* (BCRV) and the proposed species “*Bacopa chlorosis virus*” [6].

13 TSV has a broad host range that includes members of *Solanaceae* and *Leguminosae*. TSV
14 has been detected frequently amongst weeds bordering agricultural fields from where it
15 serves as a source of inoculum that could infect crops. Thrips *Microcephalothrips*
16 *abdominalis* and *Thrips tabaci* are thought to play roles in the transmission of TSV [10].
17 The mechanism of transmission is described mainly as mechanical in nature, where TSV-
18 contaminated pollen grains are introduced into the plant through wounds [11-13].

19 In the diagnosis of novel, unidentified or unusual viral plant diseases, causative agents
20 may go undetected when using methods such as PCR or ELISA, because these tests are
21 very specific to a particular species or even strain of a virus. Also, given the fast rate of
22 evolution of many RNA viruses, methods to be developed need to detect and identify

1 viral variants in a broader sense [14]. Methods such as electron microscopy (EM) or sap
2 inoculation of plant virus indicator species do not allow accurate species level diagnosis.
3 This research explores Next-generation sequencing (NGS) which offers an alternative
4 solution where sequences are generated in a non-specific fashion and identification is
5 made based on similarity alignment against GenBank data [13,15].

6 TSV is well known to be genetically diverse [9]. This diversity could contribute to failure
7 of TSV identification by conventional methods like ELISA. In such cases, detection of
8 genetically diverse viruses can be facilitated by NGS. Rapid diagnosis can be difficult in
9 cases where a virus is not well recognized as infecting a particular plant host or
10 geographical region and is further complicated where the virus differs from known
11 strains. Currently, most methods in use for virus detection are species-specific [16] (e.g.
12 PCR or ELISA). A few detection techniques have been developed for genera and families
13 [11]. Also, sap mechanical inoculation and PCR with degenerate primers, can detect
14 viruses in a species-specific manner but are limited by the fact that they can be applied to
15 a relatively small range of potential viruses. NGS technology has proved to be a
16 powerful, sensitive technology [13,17,18] especially suited to detecting unknown or
17 unsuspected viruses.

18 1. Methodology:

19 2.1 Sample collection

20 Symptomatic leaf samples from 57 soybean plants, used in this study, were collected in a
21 soybean survey in July 2012 made in Noble and Kay counties of Oklahoma USA and
22 were negative by dot-immunobinding assay (DIBA) to 15 viruses known to infect

1 soybeans (Ali et al. 2014, in preparation). Soybean plants were at the young stage before
2 flowering. Field symptoms of a typical virus-infected plant included mild mottling,
3 rugosity and chlorotic lesions on the third tri-foliolate compound leaf. However, all
4 symptomatic plants tested negative for TSV by DIBA. To further investigate the cause of
5 the symptoms and whether an unknown virus was causing the disease symptoms, NGS
6 was performed.

7 2.2 454-sequencing

8 Total RNA was extracted from leaf samples using Tri Reagent (Molecular Research
9 Center, Cincinnati, OH). The RNA concentrations of the samples were estimated by
10 Nanodrop (Thermo scientific, Waltham, MA). Nanodrop concentrations ranged from
11 600ng/μl to 1700ng/μl. Four samples of RNA concentrations greater than 1 μg/μl were
12 pooled for sequencing. These samples were chosen because of their relatively high RNA
13 yields (> 1 μg/μl) and purity as evidenced by bright bands on gel electrophoretic
14 separation on a 1.5% agarose gel.

15 The pooled RNA extract mixtures were used to synthesize double-stranded cDNA using
16 the cDNA Synthesis System kit (Roche, Indianapolis, IN) following the manufacturer's
17 protocols. A library was then produced from the RNA extract, using kits supplied by 454-
18 Roche. Sequencing was performed following the manufacturer's protocols.

19 2.3 Sequence analysis

20 The sequence reads obtained were assembled using NEWBLER v. 2.7 (Roche). The
21 resulting contigs and unassembled sequences were aligned to the NCBI nr GenBank

1 viruses database (taxid: 10239) using BLAST 2.7. Three contigs (00539, 00542 and
2 00545) were chosen for further sequence and phylogenetic analyses.

3 Open reading frame prediction and protein translation analyses were performed with
4 ExPASy, a translating bioinformatics tool [19]. Sequence alignments and phylogenetic
5 trees were constructed (using a maximum-likelihood method with 500 bootstrap
6 replicates) using Clustal W and MEGA (v. 6) [20]. Phylogenetic trees were constructed
7 using the amino acid sequences encoded by RNA 1, RNA 2 and RNA 3 of TSV isolates
8 (Table 1) and by contigs 00539, 00542 and 00545. Methods of phylogenetic analyses
9 used were neighbor-joining, maximum likelihood and maximum parsimony.
10 Evolutionary distances (number of base substitutions per site) were computed using the
11 Maximum Composite Likelihood method and were used to determine branch lengths on
12 trees. Codon positions included were 1st + 2nd + 3rd + noncoding. Initial tree(s) for the
13 heuristic searches were obtained automatically by maximum parsimony when the number
14 of common sites was below 100 or less than one-fourth of the total number of sites, or by
15 the BIONJ method with a MCL distance matrix.

16 Estimates of evolutionary divergence between sequences was carried out based on the
17 pair-wise amino acid distances per site between contigs 00539 and isolates of TSV
18 (Table1) that had produced BLAST 2.7 hits with high statistical significance (E-value
19 cutoff 10^{-3}). Standard error estimate(s) were obtained by a bootstrap procedure (500
20 replicates). Analyses were conducted using the Dayhoff matrix based model [21,20].

21 Recombination analysis was made by RDP [22] with nucleotide sequence alignments of
22 contigs 00539, 00542 and 00545 separately with the corresponding segments of other

1 TSV isolates (Table 1). When a recombination event was supported by at least three
2 different algorithms out of the seven incorporated in the RDP software the recombination
3 event was considered to be significant. These analyses were performed using the default
4 settings from the different detection programs.

5 2.4 Primer design

6 Because of inadequacies of published primers, specific PCR primers were designed
7 targeting the TSV coat protein gene: the sense primer TSV1789Fnd 5'-
8 GCTATCGTCTGCAGCCTCGA-3' (1758-1777) and the antisense primer TSV1982Rnd,
9 5'-CCACATCGCACACAGGAATT-3' (1932-1951). The primers were designed based
10 on the relatively conserved region (nt 1700-1970) of ilarvirus coat protein genes.
11 Published degenerate primers TSV1789F and TSV1982R [11] were used as reference to
12 design strain-specific primers. Web-interface applications Primer3, mFold, and BLASTn
13 were used to predict thermodynamic parameters. Internal secondary structures and self-
14 dimers were determined using mFold. Specificity was confirmed in silico by aligning the
15 primer sequences using ClustalW to contig 00545. Primers were synthesized by IDT
16 (Integrated DNA Technologies, Inc., Coralville, IA).

17 2.5 cDNA synthesis

18 First-strand cDNA synthesis was performed, using random hexamer primers on RNA
19 extracted from symptomatic soybean leaves. SuperscriptIII reverse transcriptase
20 (Invitrogen, Carlsbad, CA) was used according to the manufacturer's instructions for
21 reverse transcription of 1 µg of RNA template to synthesize the first-strand cDNA.

2.6 Endpoint RT-PCR

Endpoint RT-PCR assays were carried out with reaction mixtures that contained 10 µl of GoTaq Green Master Mix (Promega), 1 µl (5 µM) each of primers TSV1789Fnd and TSV1982Rnd, 2 µl cDNA template, and 6 µl of nuclease-free water (Ambion, Austin, TX). Reactions were performed in an MJ Research thermal cycler (PTC-100, Ramsey, MN). The cycling parameters consisted of an initial denaturation at 95°C for 1 min, followed by 35 cycles of denaturation at 95°C for 20 s, annealing at 62°C for 20 s and 56°C for 10 s, extension at 72°C for 45 s and a final extension for 3 min at 72°C. Negative (non-template control; water) controls were included in each round of PCR amplification. Five microliters of the reaction volume was electrophoresed in a 1.5% agarose gel in 1X TAE buffer. The product size (193 bp) was determined using 1 Kb Plus DNA Ladder (Invitrogen). Bands were excised from the gel and the Qiaquick Gel Extraction Kit was used to recover the PCR product. The purified DNA fragments were sequenced from both directions by automated DNA sequencing using "BigDye™"-terminated reactions and Applied Biosystems 3730 Genetic Analyzer at the Recombinant DNA/Protein Resource Facility at Oklahoma State University.

2.7 Inclusivity and Exclusivity panel

Reference positive controls for *Soybean mosaic virus* (SMV), *Bean pod mottle virus* (BPMV), *Cucumber mosaic virus* (CMV) and *Tobacco mosaic virus* (TMV), (Agdia, Elkhart, IN) constituted the exclusivity panel for specificity assays. In the inclusivity panel, TSV samples from Agdia (Elkhart, IN) were used.

2.8 Plasmid positive control pMD-01

1 A plasmid positive control pMD-01 was generated for the new TSV-OK isolate by
2 cloning a partial gene sequence (193 bp) of TSV coat protein amplified from first strand
3 cDNA using primer set TSV1789Fnd and TSV1982Rnd. Resulting amplicons were
4 excised and eluted from the agarose gel using a Qiagen gel extraction kit and cloned into
5 a plasmid vector (pCR2.1-TOPO) using a TOPO-TA Cloning Kit (Invitrogen, Grand
6 Island, NY). Plasmid DNA pMD-01 carrying target gene sequence was purified from
7 overnight grown bacterial cultures using a QIAprep Spin Miniprep Kit (Qiagen) and was
8 sequenced at the Oklahoma State University Recombinant DNA/Protein Resource
9 Facility using primers TSV1789Fnd and TSV1982Rnd. The resultant sequence aligned
10 perfectly (100% nucleotide identity) with the target gene and plasmid sequences.

11 2.9 Quantification of target cDNA with SYBRgreen RT-qPCR

12 The concentrations of plasmid pMD-01 were measured using a NanoDrop v.2000
13 spectrophotometer (Thermo Fisher Scientific, Inc., Worcester, MA). One ng/ μ l of pMD-
14 01 was used to quantify the target pathogen genomic nucleic acid with RT-qPCR, using
15 cDNA synthesized from infected plant RNA. The plasmid equivalent concentration
16 obtained from the C_T value was multiplied by a factor of 2 to yield concentration of TSV
17 viral genomes in the sample. RT-qPCR amplification was carried out in 20 μ l reaction
18 mixtures containing 10 μ l of Platinum SYBR Green, qPCR SuperMix-UDG (Invitrogen),
19 0.5 μ l (10 μ M) each of primers TSV1789Fnd and TSV1982Rnd, 2 μ l of template cDNA,
20 and 6.4 μ l of nuclease-free water. Positive and negative controls were included in each
21 round of PCR amplification, and each reaction was performed in three replicates. The
22 cycling parameters included two initial holds each for 2 min at 50 and 95°C, followed by
23 40 cycles of 95°C for 15 s and 60°C for 45 s.

1 2.10 Sensitivity assays

2 To determine the assay detection limits using primer set TSV1789Fnd and TSV1982Rnd,
3 two sensitivity assays were performed. Target cDNA from TSV-infected soybean leaves
4 made using TSV1982Rnd was quantified as described above and serially diluted in 10-
5 fold increments, from 1 ng to 1 fg per reaction. A standard curve was generated using
6 target TSV-OK cDNA and plasmid pMD-01 (plasmid pMD-01 also serially diluted in 10-
7 fold increments and used at 1 ng to 1 fg per reaction). Each RT-qPCR was performed in
8 three replicates. Assay sensitivity was also tested using primers and methods described
9 in Pappu et al. [23] since these primers were successful in detecting TSV-OK.

10 3. Results:

11 3.1 Analysis of total nucleic acids and 454 sequencing

12 Pooled RNA from symptomatic, but DIBA negative, samples were sequenced by 454-
13 sequencing and 238,946 reads were generated. The average length of reads was 480bp.
14 The reads were assembled into 600 contigs and the raw reads and the contigs were
15 aligned with the NCBI GenBank nr database using BLAST 2.7. Contigs 00539
16 (length=3431), 00542 (length=2857) and 00545 (length=2162) were found to have
17 significant similarity to sequences of TSV RNA 1, RNA 2 and RNA 3, respectively. Hits
18 to isolates of TSV provided 99% coverage, 89% identity and an E value of 0.0.
19 Strawberry necrotic shock virus and other ilarvirus member sequences were among the
20 hits detected but these had low query coverage (approximately 40%) and identity
21 (approximately 60%). Other hits with viral genomes had poor statistical significance

1 based on E value, % identity and query coverage. Using the three contigs as queries in a
2 BLASTn search of sequences accumulated during the Plant Virus Biodiversity and
3 Ecology project [24] revealed, from a single plant specimen of *Coreopsis grandiflora*
4 collected in 2005 from The Tallgrass Prairie Preserve (TPP) of Osage Co. Oklahoma,
5 adjacent to Kay County, four sequence reads covering residues 459-579 of contig 00545
6 with 98.2% identity (Supplemental file), whereas identity values compared with RNA 3
7 sequences from other TSV strains were between 81 and 83%.

8 3.2 Phylogenetic analyses

9 The resulting hits of complete genomic segments from the BLAST analysis (E value cut
10 off 10^{-3}) of contigs 00539, 00542 and 00545, were further analyzed. An analysis of the
11 pair-wise amino acid distances between replicase sequences of pairs of TSV isolates
12 (Table 2) showed that contig 00539 and TSV 1973 are more distant from other isolates
13 than those isolates are from one another. Similar results were obtained with RdRp and CP
14 sequences encoded by contigs 00542 and 00545 respectively (Tables S2 and S3). MEGA
15 6 analyses were based on the amino acid sequences of the polypeptides encoded by RNA
16 1 and RNA 2 and the CP ORF of RNA 3. In the case of RNA 1 (replicase), it was found
17 that all isolates from India clustered together (Fig. 1A). Isolates from the USA (Illinois
18 and Henry) clustered with TSV 2334 from eastern Australia. Contig 00539 and another
19 eastern Australian isolate, TSV isolate 1973, formed a distinct clade, setting these two
20 isolates apart from the rest of the TSV isolates. A similar pattern was seen in the trees
21 built with the RdRp sequence of contig 00542 polypeptide (Figure S1) and with the CP of
22 contig 00545 (Figure 1B).

3.3 Recombination analysis

For RNA 1, two separate recombination events were detected in contig 00539 (Table 3). One of these events was seen in the 5'UTR region. The major and minor parents were TSV Illinois and TSV 1973 from eastern Australia. Two other events were also detected in the RNA 2. One event was in TSV pumpkin from India (RdRp, 387-528) and the parental isolates involved were TSV 2334 (major parent) and contig 00542 (minor parent). The second recombination event was detected in TSV 2334 (RdRp, 2914-2978) and the isolates involved were TSV pumpkin and TSV Henry (Kentucky). Three other separate recombination events were detected in RNA 3. One event is located in the CP coding region of RNA 3 (nt 1513-2012). This recombination was found in isolates TSV Henry and TSV Illinois. The major and minor parents involved were TSV India and TSV (N), respectively. The second recombination event is located in the MP coding region of RNA 3 (nt 369-604). This recombination event was present in four Indian isolates: TSV Bangalore, TSV Okra, TSVAPGA and TSV India. The major and minor parents are TSV Henry and TSV (N). The third recombination was observed in the CP coding region of RNA 3 (nt 1294-1591). The major and minor parents involved were Contig 00545 and TSV 2334, respectively. All these instances of recombination are supported by at least 3 methods (Table 3) employed by RDP3 in detecting recombination and met the criteria to be considered true recombinants. The average P-values of each event for each detection program are shown in Table 3.

3.4 Detection using End point PCR and SYBR green qRT-PCR

1 Sequences of published TSV primer pairs previously used for detection [25,23,11,26,27]
2 were aligned to contigs 00539, 00542 and 00545, respectively for RNA 1, RNA 2 and
3 RNA 3. No primers matched perfectly. Each primer for contigs 00539 and 00542 had two
4 or more mismatched positions with published primers. Hence, these contigs were not
5 considered as targets for amplification. For contig 00545, there was one instance where
6 there was just one mismatch and two instances where there were two mismatches (Table
7 S1). Considering that alignments of published primers, with contig 00545 exhibited
8 greater amplification potential than with contigs 00539 and 00542, contig 00545 was
9 selected as the target gene for primer design.

10 Previously designed primers with less than three mismatches were used for amplification.
11 Endpoint RT-PCR was carried out using primers and methods described in Ravi et al.
12 [25], Sharman et al. [11] and Pappu et al. [23]. The degenerate primers TSV1789F and
13 TSV1982R [11] had only one mismatch in the reverse primer (Table S1) but failed to
14 amplify the product. Detection was successful only with primers and methods mentioned
15 in Pappu et al. [23]. Primers designed by Pappu et al. [23] were able to amplify the target
16 region in the cDNA and gave a 600 bp product but lacked sensitivity. Hence an endpoint
17 RT-PCR method was developed with newly-designed strain-specific primers
18 TSV1789Fnd and TSV1982Rnd.

19 The redesigned TSV specific primers TSV1789Fnd and TSV1982Rnd yielded sharp
20 bright bands at 193 bp. Bands were excised and purified PCR products were sequenced.
21 The resulting sequence had complete nucleotide identity to the target region of
22 contig00545. Twenty-three of 57 DIBA TSV-negative samples were tested using primers
23 TSV1789Fnd and TSV1982Rnd and were found to be positive for TSV by this test. The

1 samples were also tested with SYBR green qRT-PCR using primers TSV1789Fnd and
2 TSV1982Rnd and found to be positive. All RT-PCR tests with the exclusivity panel
3 tested negative and the RT-PCR test with the positive control (TSV from Agdia, Elkhart,
4 IN) was positive (Figure S3).

5 3.5 Sensitivity assays

6 cDNAs synthesized from total RNA of virus-infected samples were quantified by
7 generating a standard curve with known quantities of pMD-01. This step allows
8 standardization of assay sensitivity when performed using plant samples in which the
9 virus titer is unknown and/or possibly low. The target cDNA and plasmid pMD-01 were
10 serially diluted in 10-fold increments from 1ng to 1 fg for each reaction. End point RT-
11 PCR was able to detect virus titer as low as 1 fg in case of both cDNA and plasmid pMD-
12 01 (Figure 2). A standard curve was generated in SYBR green qRT-PCR assay using
13 target cDNA and plasmid pMD-01 serially diluted in 10 fold increments from 1ng to 1 fg.
14 Primers TSV1789Fnd and TSV1982Rnd in SYBR green qRT-PCR assays were able to
15 detect the target DNA up to 1 fg with both plasmid pMD-01 and cDNA (Figure 3). Each
16 qRT-PCR assay was performed in three replicates. The R^2 values were 0.95 and 0.96 for
17 plasmid pMD-01 and cDNA respectively (Table S4). Plasmid pMD-01 and cDNA had
18 slopes of -2.2 and -2.1 respectively and they both had reaction efficiencies of 1.8. An end
19 point PCR was also performed with primers and methods described in Pappu et al. [23]
20 and these reactions were found to be sensitive up to 10 ng (Figure S2).

21 4. Discussion

1 Accurate diagnosis is necessary for effective management of disease. NGS can provide a
2 useful adjunct to diagnosis of viral diseases of plants for several reasons. Multiple viruses
3 transmitted by different vectors can cause similar symptoms and thus require different
4 management strategies. Often, virus infections in plants can be asymptomatic or
5 incipiently symptomatic. However, the potential to turn into severe outbreaks persist by
6 developing obvious symptoms given the right factors and conditions (tolerance level of
7 the plant, inability of the virus to adapt to the host, presence of vectors and other
8 environmental factors). Virus evolution, especially of RNA viruses, is driven by
9 mutations, recombination and reassortment [28] and can influence the molecular
10 characteristics used as basis for their detection by methods like RT-PCR and ELISA.
11 General antibodies can fail to detect some strains of a particular virus due to changes in
12 antigenic determinants. In RT-PCR if the primers are not a good match to the template
13 they may fail to amplify the target gene or, when the primers are not sensitive enough for
14 diagnostic purposes when the viral titer is low. In these instances, NGS is an alternative
15 rapid tool for accurate diagnosis.

16 Contigs 00539, 00542 and 00545 define a molecularly distinct TSV strain as shown by
17 the pair-wise amino acid distances per site, phylogenetic and recombination analyses. In
18 the phylogenetic tree of RNA 1 (Fig. 1), TSV 1973 and contig 00539 share a branch
19 connecting them to the rest of the tree, even though a greater similarity between them
20 than to other isolates is not evident in their pair-wise amino acid distance analysis (Table
21 2). This discrepancy could be the result of the presence of a recombination event between
22 these two strains in the 5' UTR region and the RNA 1 helicase region (Table 3). We also
23 detected the presence of recombination in the CP region between these two strains. In the

1 CP phylogenetic tree (Fig. 1), TSV 1973 and contig 00545 each form a long branch, and
2 share a clade connecting them to the rest of the tree, even though their pair-wise amino
3 acid distances per site do not indicate that TSV 1973 and contig 00545 have a greater
4 similarity to one another than to other isolates (Table S3).

5 The fact that the TSV strain discovered in Oklahoma's soybean crop (or a close relative)
6 was also found in a nature preserve, the Tallgrass Prairie Preserve, on which soybeans
7 had never been grown suggests that the Oklahoma disease may arise from movement of
8 the virus from wild plants into cultivated soybean fields or from the soybean plants to the
9 Preserve plants [24]. Further biological characterization of TSV-OK is in progress.

10 While NGS allows detection of most organisms present in a sample, it suffers from time-
11 consuming limitations during assembly and analysis of sequence data suggesting that
12 NGS based virus discovery will only become more complicated with time. To circumvent
13 these limitations, it has been shown that unique pathogen-specific sequences can be used
14 in searches of unassembled, non-quality checked, sequence data to accurately detect
15 pathogens in plant samples [29]. This process, termed E-probe Diagnostic Nucleic acid
16 Analysis, treats the NGS sample sequences as the "reference" database that are queried
17 with pathogen specific sequences termed E-probes. This process was very efficient at
18 detecting known RNA and DNA viruses from infected plants [30]. However, the process
19 only detects viruses for which the user directly queries using virus specific e-probes. We
20 are testing the use of broad detection E-probes developed from probes designed for a
21 virus microarray study (manuscript in preparation). These are general probes that allow
22 detection of viruses at the family level. Preliminary results show efficient and accurate
23 first line diagnostic testing can be done following this procedure.

Acknowledgements

This study was funded by the USDA-CSREES Plant Biosecurity Program (grant number 2010-85605-20542) and the Oklahoma Agricultural Experiment Station. H. Hwang and the Bioinformatics Core Facility are acknowledged for providing exceptional consistency and skill while preparing libraries and performing pyrosequencing. The authors thank the National Science Foundation (NSF-DEB 842703) and the Robert J. Sirny professorship for additional support. The authors also thank Kay Scheets and Francisco Ochoa-Corona for internal review.

References

1. Weiden MD, Ginsberg HS (1994) Deletion of the E4 region of the genome produces adenovirus DNA concatemers. *Proc Natl Acad Sci Usa* 91 (1):153-157
2. Fagbenle HH, Ford RE (1970) Tobacco-streak-virus isolated from soybeans, Glycine-max. *Phytopathology* 60 (5):814-&
3. Ghanekar AM, Schwenk FW (1974) Seed transmission and distribution of tobacco streak virus in 6 cultivars of soybeans. *Phytopathology* 64 (1):112-114
4. Rabedaux PF, Gaska JM, Kurtzweil NC, Grau CR (2005) Seasonal progression and agronomic impact of Tobacco streak virus on soybean in Wisconsin. *Plant Disease* 89 (4):391-396. doi:10.1094/pd-89-0391

- 1 5. Wang Y, Hobbs HA, Hill CB, Domier LL, Hartman GL, Nelson RL (2005) Evaluation
2 of ancestral lines of US soybean cultivars for resistance to four soybean viruses. *Crop*
3 *Science* 45 (2):639-644
- 4 6. Pallas V, Aparicio F, Herranz MC, Sanchez-Navarro JA, Scott SW (2013) The
5 Molecular Biology of Iilarviruses. In: Maramorosch K, Murphy FA (eds) *Advances in*
6 *Virus Research*, Vol 87, vol 87. *Advances in Virus Research*. Elsevier Academic Press
7 Inc, San Diego, pp 139-181. doi:10.1016/b978-0-12-407698-3.00005-3
- 8 7. Bujarski J, Figlerowicz M, Gallitlli D, Roossinck MJ, Scott SW (2012) Bromoviridae.
9 In: King AMQ, Adams MJ, Carstens EB, Lefkowitz EJ (eds) *Virus Taxonomy*, Ninth
10 Report of the International Committee on Taxonomy of Viruses. Elsevier, Amsterdam, pp
11 965-976
- 12 8. Halk EL, Franke J (1983) Identification of serological types of Apple mosaic, *Prunus*
13 necrotic ringspot and Tobacco streak viruses with monoclonal-antibodies.
14 *Phytopathology* 73 (5):789-789
- 15 9. Clement DL, Converse RH (1986) Serological relationships among 4 tobacco streak
16 virus isolates. *Phytopathology* 76 (8):842-842
- 17 10. Jones DR (2005) Plant viruses transmitted by thrips. *European Journal of Plant*
18 *Pathology* 113 (2):119-157. doi:10.1007/s10658-005-2334-1
- 19 11. Sharman M, Thomas JE (2013) Genetic diversity of subgroup 1 ilarviruses from
20 eastern Australia. *Archives of Virology* 158 (8):1637-1647. doi:10.1007/s00705-013-
21 1628-4

12. Klose MJ, Sdoodee R, Teakle DS, Milne JR, Greber RS, Walter GH (1996) Transmission of three strains of tobacco streak ilarvirus by different thrips species using virus-infected pollen. *Journal of Phytopathology-Phytopathologische Zeitschrift* 144 (6):281-284. doi:10.1111/j.1439-0434.1996.tb01530.x
13. Adams IP, Glover RH, Monger WA, Mumford R, Jackeviciene E, Navalinskiene M, Samuitiene M, Boonham N (2009) Next-generation sequencing and metagenomic analysis: a universal diagnostic tool in plant virology. *Molecular Plant Pathology* 10 (4):537-545. doi:10.1111/j.1364-3703.2009.00545.x
14. Zheng L, Wayper PJ, Gibbs AJ, Fourment M, Rodoni BC, Gibbs MJ (2008) Accumulating variation at conserved sites in potyvirus genomes is driven by species discovery and affects degenerate primer design. *PLoS ONE* 3 (2):e1586. doi:10.1371/journal.pone.0001586
15. Ho T, Tzanetakis IE (2014) Development of a virus detection and discovery pipeline using next generation sequencing. *Virology* 471–473 (0):54-60. doi:http://dx.doi.org/10.1016/j.virol.2014.09.019
16. Arif M, Aguilar-Moreno GS, Wayadande A, Fletcher J, Ochoa-Corona FM (2014) Primer Modification Improves Rapid and Sensitive In Vitro and Field-Deployable Assays for Detection of High Plains Virus Variants. *Applied and Environmental Microbiology* 80 (1):320-327. doi:10.1128/aem.02340-13
17. Al Rwahnih M, Daubert S, Golino D, Rowhani A (2009) Deep sequencing analysis of RNAs from a grapevine showing Syrah decline symptoms reveals a multiple virus

1 infection that includes a novel virus. *Virology* 387 (2):395-401.
2 doi:10.1016/j.virol.2009.02.028

3 18. Kreuze JF, Perez A, Untiveros M, Quispe D, Fuentes S, Barker I, Simon R (2009)
4 Complete viral genome sequence and discovery of novel viruses by deep sequencing of
5 small RNAs: a generic method for diagnosis, discovery and sequencing of viruses.
6 *Virology* 388 (1):1-7

7 19. Gasteiger E, Gattiker A, Hoogland C, Ivanyi I, Appel RD, Bairoch A (2003)
8 ExPASy: The proteomics server for in-depth protein knowledge and analysis. *Nucleic
9 Acids Res* 31 (13):3784-3788

10 20. Tamura K, Stecher G, Peterson D, Filipski A, Kumar S (2013) MEGA6: Molecular
11 Evolutionary Genetics Analysis Version 6.0. *Molecular Biology and Evolution* 30
12 (12):2725-2729. doi:10.1093/molbev/mst197

13 21. Schwarz R. DM (1979) Matrices for detecting distant relationships. *Atlas of protein
14 sequences*: 353-358

15 22. Martin DP, Lemey P, Lott M, Moulton V, Posada D, Lefevre P (2010) RDP3: a
16 flexible and fast computer program for analyzing recombination. *Bioinformatics* 26
17 (19):2462-2463. doi:10.1093/bioinformatics/btq467

18 23. Pappu HR, Hammett KRW, Druffel KL (2008) Dahlia mosaic virus and Tobacco
19 streak virus in dahlia (*Dahlia variabilis*) in New Zealand. *Plant Disease* 92 (7):1138-1138.
20 doi:10.1094/pdis-92-7-1138b

- 1 24. Wren JD, Roossinck MJ, Nelson RS, Scheets K, Palmer MW, Melcher U (2006)
2 Plant virus biodiversity and ecology. PLoS Biol 4 (3):e80.
3 doi:10.1371/journal.pbio.0040080
- 4 25. Kumar NA, Narasu ML, Zehr UB, Ravi KS (2008) Molecular characterization of
5 Tobacco streak virus causing soybean necrosis in India. Indian Journal of Biotechnology
6 7 (2):214-217
- 7 26. Bhat AI, Jain RK, Ramiah M (2002) Detection of Tobacco streak virus from
8 sunflower and other crops by reverse transcription polymerase chain reaction. Indian
9 Phytopathology 55 (2):216-218
- 10 27. Hosseini S, Habibi MK, Mosahebi G, Motamedi M, Winter S (2012) First report on
11 the occurrence of Tobacco streak virus in sunflower in Iran. Journal of Plant Pathology
12 94 (3):585-589
- 13 28. Ali A, Roossinck MJ (2010) Genetic bottlenecks during systemic movement of
14 Cucumber mosaic virus vary in different host plants. Virology 404 (2):279-283.
15 doi:http://dx.doi.org/10.1016/j.virol.2010.05.017
- 16 29. Stobbe AH, Schneider WL, Hoyt PR, Melcher U (2014) Screening metagenomic data
17 for viruses using the e-probe diagnostic nucleic acid Assay (EDNA). *Phytopathology*,
18 **104**, 1125-1129
- 19 30. Stobbe AH, Daniels J, Espindola AS, Verma R, Melcher U, Ochoa-Corona F, Garzon
20 C, Fletcher J, Schneider W (2013) E-probe Diagnostic Nucleic acid Analysis (EDNA): A

- 1 theoretical approach for handling of next generation sequencing data for diagnostics.
- 2 Journal of Microbiological Methods 94 (3):356-366. doi:10.1016/j.mimet.2013.07.002
- 3

- 1 **Table 1:** Isolate identity, geographical locations and accession numbers of the RNAs
- 2 used for analyses in this article.

| TSV isolate | Geographical location | Accession numbers | | |
|--|-------------------------|-------------------|------------|------------|
| | | RNA 1 | RNA2 | RNA3 |
| Tobacco streak virus OK | Oklahoma, USA | KP256522 | KP256521 | KP256520 |
| Tobacco streak virus isolate 1973 | Eastern Australia | JX463334.1 | JX463335.1 | JX463336.1 |
| Tobacco streak virus isolate pumpkin | Karnataka, India | FJ561299.1 | FJ561300.1 | |
| Tobacco streak virus isolate okra | Karnataka, India | FJ561302.1 | | FJ561304.1 |
| Tobacco streak virus isolate REPDHP | Tamil Nadu, India | KF264473.1 | | |
| Tobacco streak virus isolate REPKAR | Tamil Nadu, South India | KF264474.1 | | |
| Tobacco streak virus isolate REPCBE2 | Tamil Nadu, South India | KF264472.1 | | |
| Tobacco streak virus isolate Henry | Kentucky, USA | JX073656.1 | JX073657.1 | JX073658.1 |
| Tobacco streak virus isolate 2334 | Eastern Australia | JX463337.1 | JX463338.1 | JX463339.1 |
| Tobacco streak virus isolate Illinois | Illinois, USA | FJ403375.1 | FJ403376.1 | FJ403377.1 |
| Tobacco streak virus from India | Karnataka, India | | | DQ067449.1 |
| Tobacco streak virus isolate TSVAPGA | Andhra Pradesh, India | | | FJ417082.1 |
| Tobacco streak virus | Karnataka, India | | | FJ655171.1 |

| | | | | |
|--|-------------------|--|--|------------|
| Tobacco streak virus | Karnataka, India | | | FJ655170.1 |
| Tobacco streak virus (N) | Netherlands | | | X00435.1 |
| Tobacco streak virus isolate 1025 | Eastern Australia | | | JX463347.1 |

1

2

Table 2: Pair-wise amino acid distances in RNA 1 of selected TSV strains. The amino acid distances per site between sequences is shown in the lower left triangle with standard error estimates in the upper right triangle. Analyses were conducted using the Dayhoff matrix based model with 500 bootstrap replicates.

| | Contig 00539 | TSV Illinois | TSV 2334 | TSV Henry | TSV REPCBE 2 | TSV REPKAR | TSV REPDHP | TSV Okra | TSV Pumpkin | TSV 1973 |
|------------------------|-------------------------|-------------------------|---------------------|----------------------|-----------------------------|-----------------------|-----------------------|---------------------|------------------------|---------------------|
| Contig00539 | | 0.008 | 0.009 | 0.009 | 0.009 | 0.009 | 0.009 | 0.009 | 0.010 | 0.011 |
| TSV Illinois | 0.082 | | 0.006 | 0.006 | 0.007 | 0.007 | 0.007 | 0.007 | 0.008 | 0.010 |
| TSV 2334 | 0.083 | 0.042 | | 0.006 | 0.007 | 0.007 | 0.007 | 0.007 | 0.007 | 0.010 |
| TSV Henry | 0.083 | 0.043 | 0.045 | | 0.008 | 0.008 | 0.008 | 0.008 | 0.008 | 0.011 |
| TSV REPCBE2 | 0.092 | 0.068 | 0.061 | 0.069 | | 0.001 | 0.001 | 0.002 | 0.003 | 0.011 |
| TSV REPKAR | 0.093 | 0.069 | 0.062 | 0.070 | 0.001 | | 0.002 | 0.002 | 0.003 | 0.011 |
| TSV REPDHP | 0.093 | 0.068 | 0.062 | 0.070 | 0.003 | 0.003 | | 0.002 | 0.003 | 0.011 |
| TSV Okra | 0.094 | 0.069 | 0.063 | 0.071 | 0.005 | 0.006 | 0.006 | | 0.004 | 0.011 |
| TSV Pumpkin | 0.099 | 0.074 | 0.068 | 0.075 | 0.010 | 0.010 | 0.010 | 0.013 | | 0.011 |
| TSV 1973 | 0.126 | 0.120 | 0.119 | 0.123 | 0.139 | 0.140 | 0.140 | 0.141 | 0.146 | |

Table 3: Analysis of recombination of RNA sequences of TSV isolates

| Virus isolate | Genomic region | Event | Major/minor parents | RDP analysis | Average P- value |
|---|-----------------------|--------------|----------------------------|---|--|
| Consensus 00539 | RNA 1 | 1 | TSV Illinois/TSV 1973 | RDP Geneconv Bootscan Maxchi Chimaera | 3.27×10^{-7} 5.56×10^{-5} 6.02×10^{-6} 1.75×10^{-3} 7.00×10^{-5} |
| Consensus 00539 | RNA 1 (3037-3168) | 1 | TSV Illinois/TSV 1973 | RDP Maxchi Chimaera | 5.43×10^{-3} 2.22×10^{-3} 6.99×10^{-3} |
| TSV pumpkin | RNA 2 (387-528) | 1 | TSV 2334/Consensus 00542 | RDP Maxchi Chimaera | 5.19×10^{-3} 2.33×10^{-3} 6.01×10^{-3} |
| TSV 2334 | RNA 2 (2914-2978) | 1 | TSV pumpkin/TSV Henry | RDP Geneconv Bootscan Maxchi Chimaera | 3.70×10^{-5} 8.89×10^{-8} 1.59×10^{-8} 1.03×10^{-4} 2.01×10^{-2} |
| TSV Henry TSV Illinois | RNA 3 (1513-2012) | 2 | TSV India/TSV (N) | RDP Bootscan Maxchi | 1.20×10^{-2} 4.43×10^{-2} 6.25×10^{-3} |
| TSV Bangalore TSV Okra TSVAPGA TSV India | RNA 3 (369-604) | 4 | TSV Henry/TSV (N) | Maxchi Chimaera SiScan | 8.78×10^{-3} 2.24×10^{-2} 2.17×10^{-2} |
| TSV 1973 | RNA 3 (1294-1591) | 1 | Consensus 00545/TSV 2334 | RDP Bootscan Maxchi Chimaera SiScan | 2.86×10^{-4} 2.28×10^{-2} 7.32×10^{-7} 2.38×10^{-5} 4.28×10^{-11} |

Figure Legends

Fig 1: Phylogenetic relationships among TSV strains: (A) Amino acid sequences encoded by contig 00539 and different isolates of TSV with statistically significant relatedness; (B) Predicted CP sequences encoded by contig 00545 and other isolates of TSV with statistically significant relatedness. Trees shown were inferred by using the maximum-likelihood method of MEGA6 based on the Tamura-Nei model. Trees with the highest log likelihood (-16662.1875) are shown. Branches corresponding to partitions reproduced in more than 50 % of the 500 bootstrap replicates are numbered. Abbreviations used are as follows: TSV (Tobacco streak virus) followed by the name of the isolate.

Fig 2: Endpoint RT-PCR sensitivity assays. Target cDNA (A) from TSV infected soybean leaves and plasmid pMD-01 (B) carrying the target gene fragment were serially diluted in 10-fold increments and used at 1 ng to 1 fg per reaction. Lane L, 1-kb ladder (Invitrogen). Lanes 1 to 7 contain serially diluted (10-fold increments) cDNA or plasmid pMD-01 at concentrations from 1 ng to 1 fg.

Fig 3: Sensitivity assays using real-time SYBR green RT-qPCR for TSV. Target cDNA (A) from TSV-infected soybean leaves and plasmid pMD-01 (B) carrying the target gene fragment was serially diluted in 10-fold increments and used at 1 ng to 1 fg per reaction. R^2 is linear correlation, M is the slope, and Ex is reaction efficiency. Each reaction was performed in three replicates. The x axes show the concentrations used, and the y axes show the observed cycle thresholds.

Figure 1A

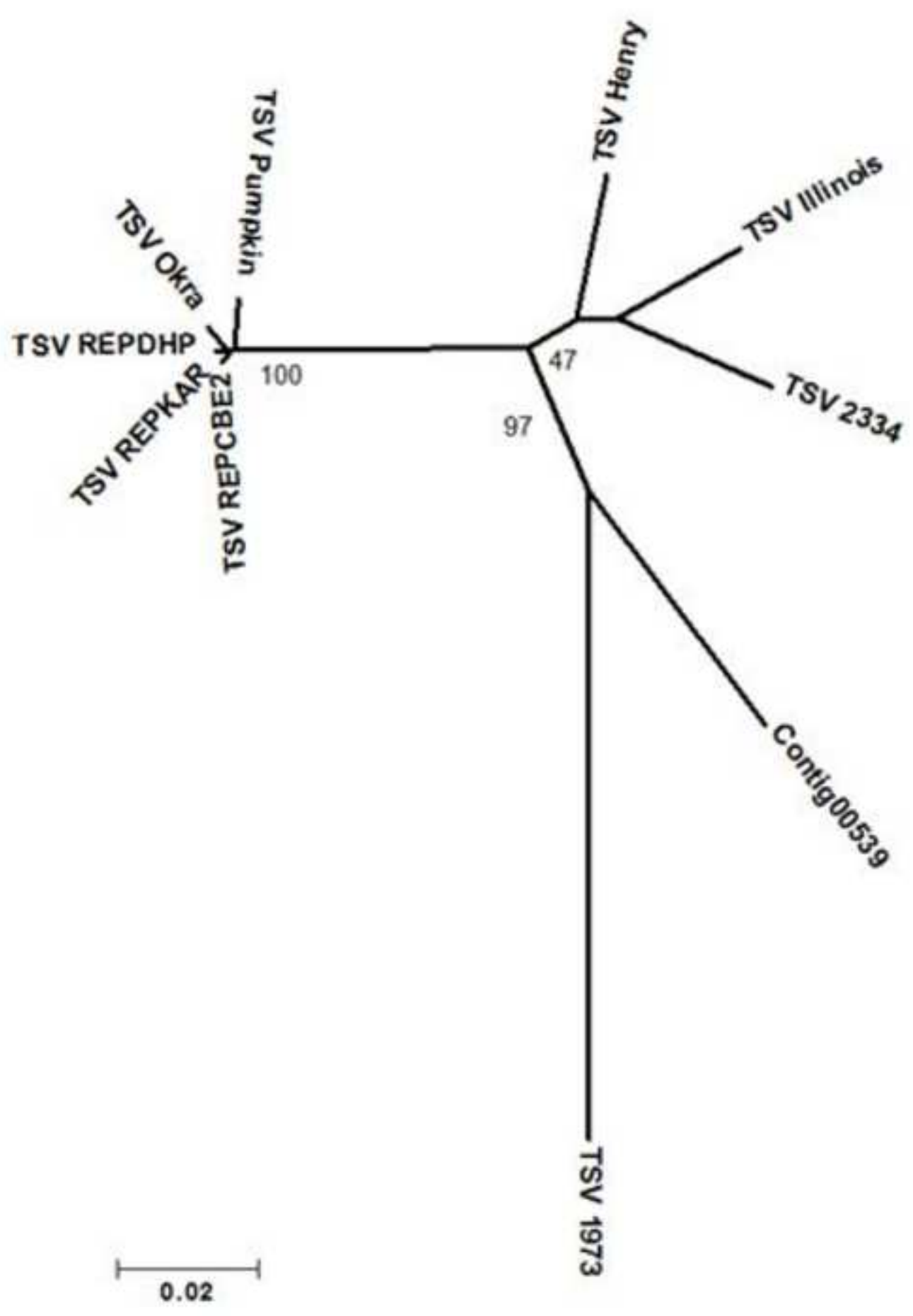


Figure 1B

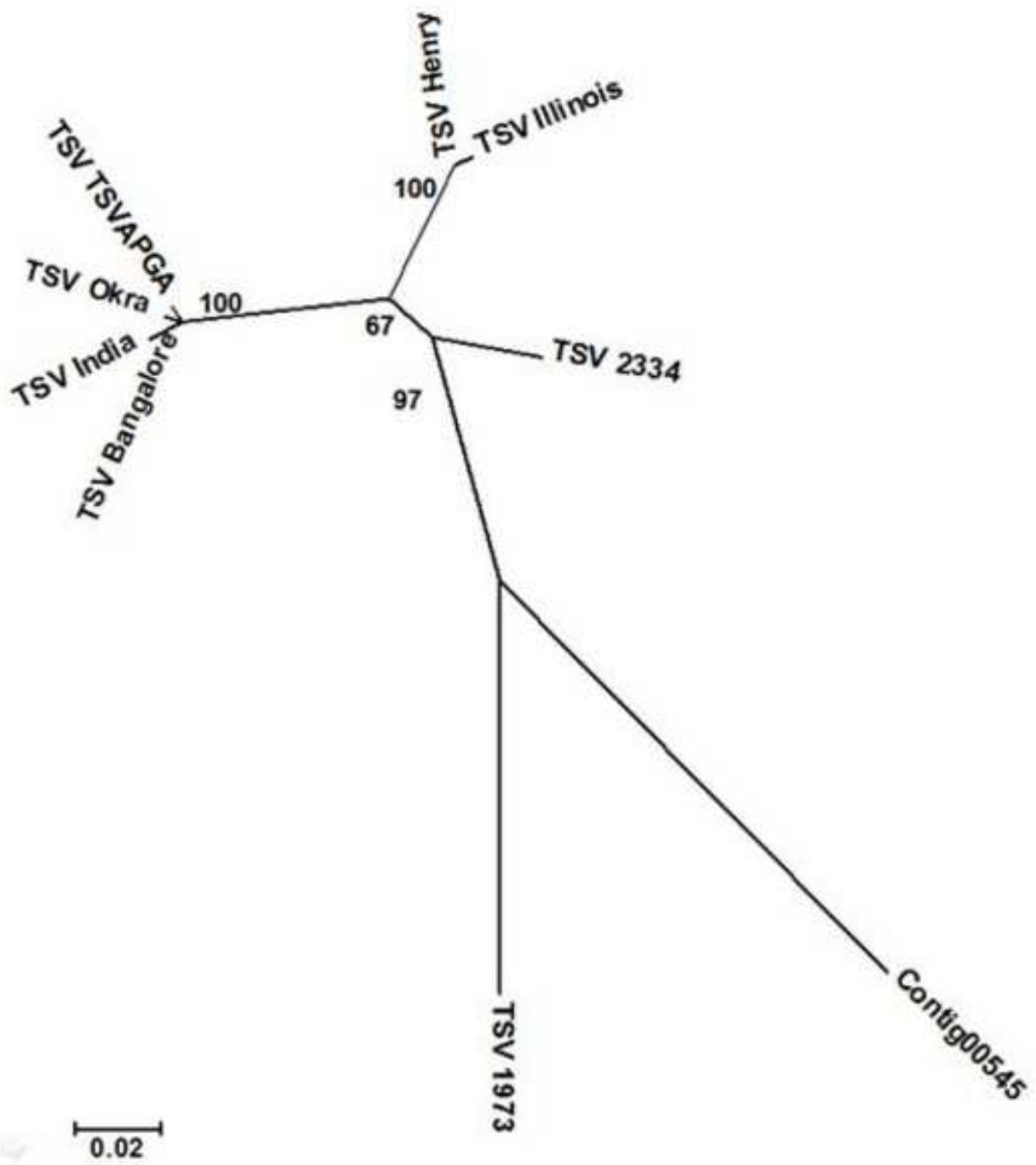


Figure 2A

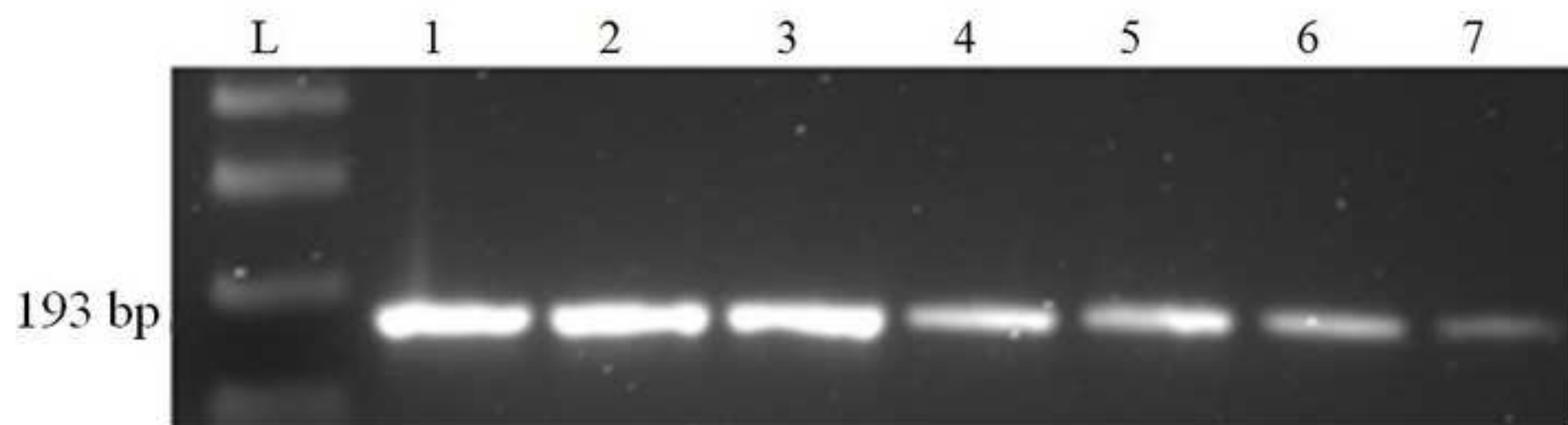


Figure 2B

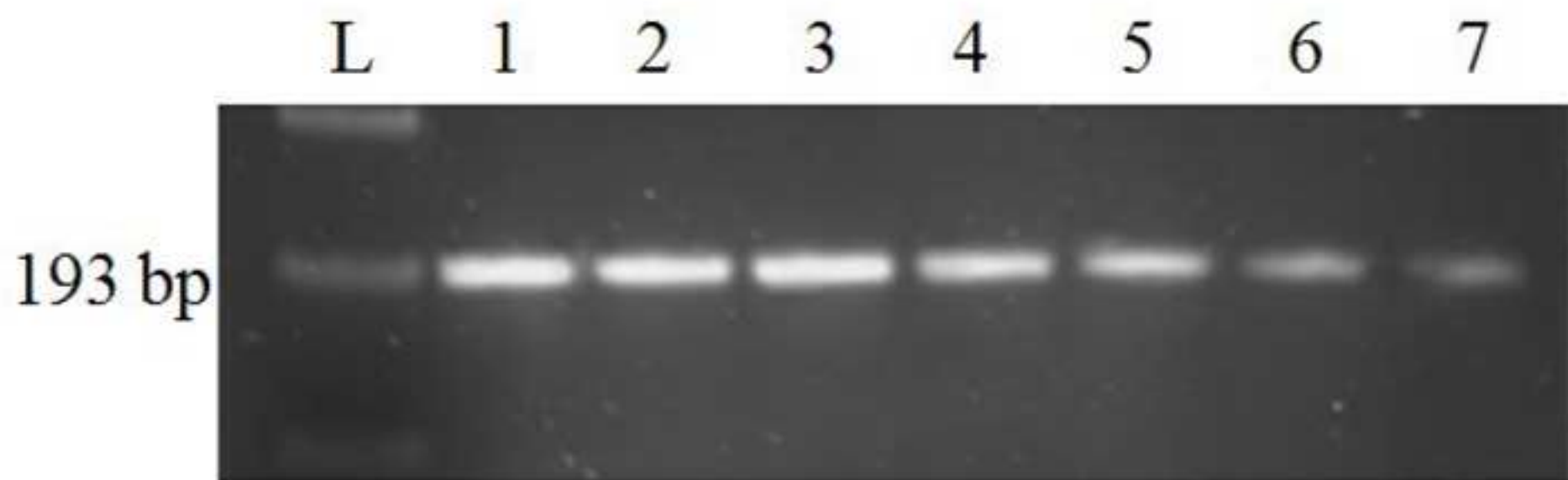


Figure 3A

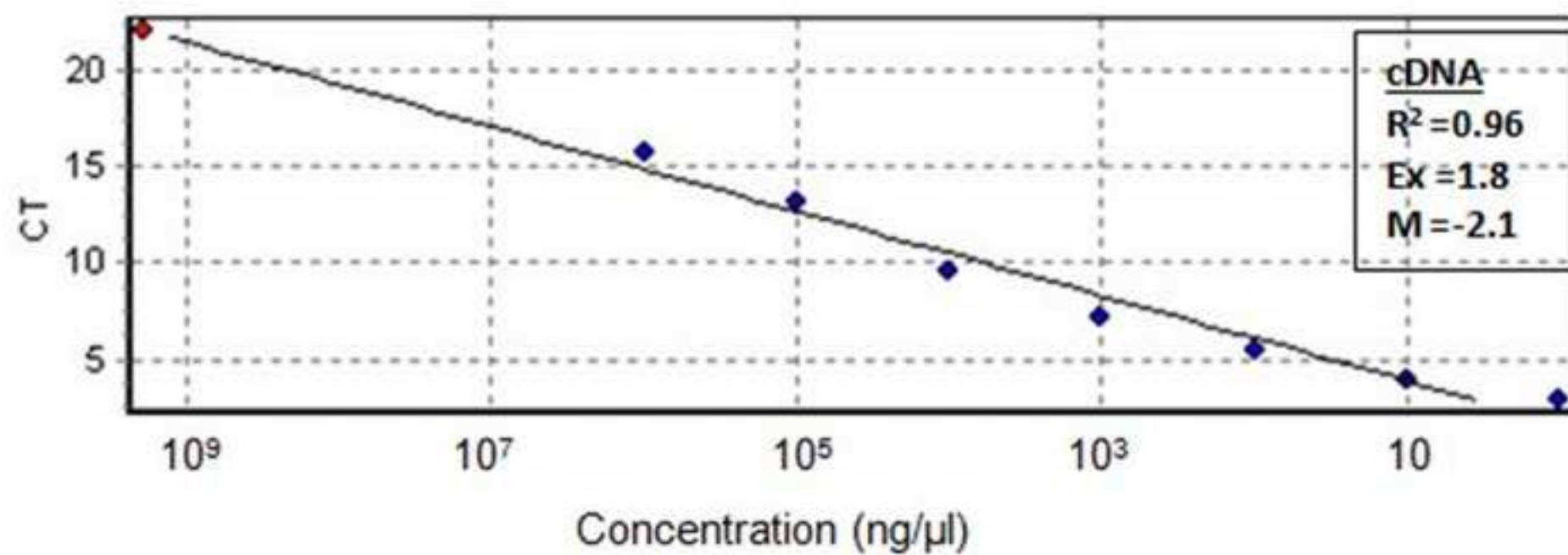


Figure 3B

